

# Scalable Approximations for Generalized Linear Problems

Murat A. Erdogdu\*      Mohsen Bayati†      Lee H. Dicker‡

November 16, 2016

## Abstract

In stochastic optimization, the population risk is generally approximated by the empirical risk. However, in the large-scale setting, minimization of the empirical risk may be computationally restrictive. In this paper, we design an efficient algorithm to approximate the population risk minimizer in generalized linear problems such as binary classification with surrogate losses and generalized linear regression models. We focus on large-scale problems, where the iterative minimization of the empirical risk is computationally intractable, i.e., the number of observations  $n$  is much larger than the dimension of the parameter  $p$ , i.e.  $n \gg p \gg 1$ . We show that under random sub-Gaussian design, the true minimizer of the population risk is approximately proportional to the corresponding ordinary least squares (OLS) estimator. Using this relation, we design an algorithm that achieves the same accuracy as the empirical risk minimizer through iterations that attain up to a cubic convergence rate, and that are cheaper than any batch optimization algorithm by at least a factor of  $\mathcal{O}(p)$ . We provide theoretical guarantees for our algorithm, and analyze the convergence behavior in terms of data dimensions. Finally, we demonstrate the performance of our algorithm on well-known classification and regression problems, through extensive numerical studies on large-scale datasets, and show that it achieves the highest performance compared to several other widely used and specialized optimization algorithms.

## 1 Introduction

We consider the following optimization problem

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad R(\beta) := \mathbb{E} [\Psi(\langle x, \beta \rangle) - y \langle x, \beta \rangle], \quad (1.1)$$

where  $\Psi : \mathbb{R} \rightarrow \mathbb{R}$  is a non-linear function,  $y \in \mathcal{Y} \subset \mathbb{R}$  denotes the response variable,  $x \in \mathcal{X} \subset \mathbb{R}^p$  denotes the predictor (or covariate), and the expectation is over the joint distribution of  $(y, x)$ . The above minimization is called a generalized linear problem in its canonical representation, and it is commonly encountered in the statistical learning. Celebrated examples include binary classification with smooth surrogate losses [BSS05, RW10], and generalized linear models (GLMs)

---

\*Department of Statistics, Stanford University, [erdogdu@stanford.edu](mailto:erdogdu@stanford.edu)

†Graduate School of Business, Stanford University, [bayati@stanford.edu](mailto:bayati@stanford.edu)

‡Department of Statistics and Biostatistics, Rutgers University and Amazon (Work conducted while at Rutgers University), [ldicker@stat.rutgers.edu](mailto:ldicker@stat.rutgers.edu),

such as Poisson regression, logistic regression, ordinary least squares, multinomial regression and many applications involving graphical models [NB72, MN89, WJ08, KF09]. These methods play a crucial role in numerous machine learning and statistics problems, and provide a miscellaneous framework for many regression and classification tasks.

The exact minimization of the stochastic optimization problem (1.1), requires the knowledge of the underlying distribution of the variables  $(y, x)$ . In practice, however, the joint distribution is not available. Therefore, after observing  $n$  independent data points  $(y_i, x_i)$ , the standard approach is to minimize the empirical risk approximation given as

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \hat{R}(\beta) := \frac{1}{n} \sum_{i=1}^n \Psi(\langle x_i, \beta \rangle) - y_i \langle x_i, \beta \rangle. \quad (1.2)$$

In the case of GLMs, the empirical risk minimization given in (1.2) is called the maximum likelihood estimation, whereas in the case of binary classification, it is generally referred to as surrogate loss minimization. Due to non-linear structure of the optimization task given in (1.2), for both problems, the minimization of the empirical risk requires iterative methods. Regardless of the problem formulation, the most commonly used optimization method is the Newton-Raphson method, which may be viewed as a reweighted least squares algorithm [MN89, BSS05]. This method uses a second-order approximation to benefit from the curvature of the log-likelihood and achieves locally quadratic convergence. A drawback of this approach is its excessive per-iteration cost of  $\mathcal{O}(np^2)$ . To remedy this, Hessian-free Krylov sub-space based methods such as conjugate gradient and minimal residual are used, but the resulting direction is imprecise [HS52, PS75, Mar10]. On the other hand, first-order approximation yields the gradient descent algorithm, which attains a linear convergence rate with  $\mathcal{O}(np)$  per-iteration cost. Although its convergence rate is slow compared to that of the second-order methods, its modest per-iteration cost makes it practical for large-scale problems. In the regime  $n \gg p$ , another popular optimization technique is the class of Quasi-Newton methods [Bis95, Nes04], which can attain a per-iteration cost of  $\mathcal{O}(np)$ , and the convergence rate is locally super-linear; a well-known member of this class of methods is the BFGS algorithm [Bro70, Fle70, Gol70, Sha70]. There are recent studies that exploit the special structure of GLMs [Erd15a], and achieve near-quadratic convergence with a per-iteration cost of  $\mathcal{O}(np)$ , and an additional cost of covariance estimation.

In this paper, we take an alternative approach for minimizing (1.1), based on an identity that is well-known in some areas of statistics, but appears to have received relatively little attention for its computational implications in large-scale problems. Let  $\beta^{\text{pop}}$  denote the true minimizer of the population risk given in (1.1), and let  $\beta^{\text{ols}}$  denote the corresponding ordinary least squares (OLS) coefficients defined as  $\beta^{\text{ols}} = \mathbb{E}[xx^T]^{-1} \mathbb{E}[xy]$ . Then, under certain random predictor (design) models,

$$\beta^{\text{pop}} \propto \beta^{\text{ols}}. \quad (1.3)$$

For logistic regression with Gaussian design (which is equivalent to Fisher's discriminant analysis), (1.3) was noted by Fisher in the 1930s [Fis36]; a more general formulation for models with Gaussian design is given in [Bri82]. The relationship (1.3) suggests that if the constant of proportionality is known, then  $\beta^{\text{pop}}$  can be estimated by computing the OLS estimator, which may be substantially simpler than minimizing the empirical risk. In fact, in some applications like binary classification, it may not be necessary to find the constant of proportionality in (1.3). Our work in this paper builds on this idea.

Our contributions can be summarized as follows.

1. We show that  $\beta^{\text{pop}}$  is approximately proportional to  $\beta^{\text{ols}}$  in the random design setting, regardless of the covariate (predictor) distribution. That is, we prove

$$\left\| \beta^{\text{pop}} - c_{\Psi} \times \beta^{\text{ols}} \right\|_{\infty} \lesssim \frac{1}{p},$$

for some  $c_{\Psi} \in \mathbb{R}$  which depends on the non-linearity  $\Psi$ . Our generalization uses zero-bias transformations [GR97]. We also show that the above relation still holds under certain types of regularization.

2. We design a computationally efficient estimator for  $\beta^{\text{pop}}$  by first estimating the OLS coefficients, and then estimating the proportionality constant  $c_{\Psi}$  via line search. We refer to the resulting estimator as the Scaled Least Squares (SLS) estimator and denote it by  $\hat{\beta}^{\text{sls}}$ . After estimating the OLS coefficients, the second step of our algorithm involves finding a root of a real valued function; this can be accomplished using iterative methods with up to a cubic convergence rate and only  $\mathcal{O}(n)$  per-iteration cost. This is cheaper than the classical batch methods mentioned above by at least a factor of  $\mathcal{O}(p)$ .
3. For random design with sub-Gaussian predictors, we show that

$$\left\| \hat{\beta}^{\text{sls}} - \beta^{\text{pop}} \right\|_{\infty} \lesssim \frac{1}{p} + \sqrt{\frac{p}{n/\log(n)}}.$$

This bound characterizes the performance of the proposed estimator in terms of data dimensions, and justifies the use of the algorithm in the regime  $n \gg p \gg 1$ .

4. We demonstrate how to transform a binary classification problem with smooth surrogate loss into a generalized linear problem, and how our methods can be applied to obtain a computationally efficient optimization scheme. We further discuss the canonicalization of the square loss, which may be of independent interest to non-convex optimization community.
5. We propose a scalable algorithm for converting one generalized linear problem to another by exploiting the proportionality relation (1.3). The proposed algorithm requires only  $\mathcal{O}(n)$  per each iteration, with no additional cost.
6. We study the statistical and computational performance of  $\hat{\beta}^{\text{sls}}$ , and compare it to that of the empirical risk minimizer (using several well-known implementations), on a variety of large-scale datasets.

The rest of the paper is organized as follows: Section 1.1 surveys the related work and Section 2 introduces the required background and the notation. In Section 3, we provide the intuition behind the relationship (1.3), which are based on exact calculations for the Gaussian design setting. In Section 4, we propose our algorithm and discuss its computational properties. Theoretical results are given in Section 5. In Section 6, we propose an algorithm to convert one GLM type to another. We discuss how a binary classification problem can be cast as a generalized linear problem in Section 7, and in Section 8 we propose a method to canonicalize the square loss. Section 9 provides a thorough comparison between the proposed algorithm and other existing methods. Finally, we conclude with a brief discussion in Section 10.

## 1.1 Related work

As mentioned in Section 1, the relationship (1.3) is well-known in several forms in statistics. Brillinger [Bri82] derived (1.3) for models with Gaussian predictors using Stein’s lemma. Li & Duan [LD89] studied model misspecification problems in statistics and derived (1.3) when the predictor distribution has linear conditional means (this is a slight generalization of Gaussian predictors). The relation (1.3) has led to various techniques for dimension reduction [Li91, LD09], and more recently, it has been studied by [PV15, TAH15] in the context of compressed sensing. It has been shown that the standard lasso estimator may be very effective when used in models where the relationship between the expected response and the signal is nonlinear, and the predictors (i.e. the design or sensing matrix) are Gaussian. A common theme for all of this previous work is that it focuses solely on settings where (1.3) holds exactly and the predictors are Gaussian (or, in the case of [LD89], very nearly Gaussian). Two key novelties of the present paper are (i) our focus on the computational benefits following from (1.3) for large scale problems with  $n \gg p \gg 1$ ; and (ii) our rigorous finite sample analysis of models with non-Gaussian predictors, where (1.3) is shown to be approximately valid. To the best of our knowledge, the present paper and its earlier version [EBD16] are the first to consider the relation (1.3) in the context of optimization.

## 2 Preliminaries and notation

We assume a random design setting, where the observed data consists of  $n$  random iid pairs  $(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)$ ;  $y_i \in \mathcal{Y} \subset \mathbb{R}$  is the response variable and  $x_i = (x_{i1}, \dots, x_{ip})^T \in \mathcal{X} \subset \mathbb{R}^p$  is the vector of predictors or covariates. We focus on problems where the minimization (1.1) is desirable, but we do not need to assume that  $(y_i, x_i)$  are actually drawn from a particular distribution or the corresponding statistical model (i.e. we allow for model misspecification).

$$\beta^{\text{pop}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \mathbb{E} [\Psi(\langle x_i, \beta \rangle) - y_i \langle x_i, \beta \rangle]. \quad (2.1)$$

While we make no assumptions on  $\Psi$  beyond smoothness, note that when the optimization problem is GLM, and  $\Psi$  is the cumulant generating function for  $y_i \mid x_i$ , then the problem reduces to the standard GLM with canonical link and regression parameters  $\beta^{\text{pop}}$  [MN89]. Examples of GLMs in this form include logistic regression with  $\Psi(w) = \log\{1 + e^w\}$ , Poisson regression with  $\Psi(w) = e^w$ , and linear regression (least squares) with  $\Psi(w) = w^2/2$ .

Our objective is to find a computationally efficient estimator for  $\beta^{\text{pop}}$ . The alternative estimator for  $\beta^{\text{pop}}$  proposed in this paper is related to the OLS coefficient vector, which is defined by  $\beta^{\text{ols}} := \mathbb{E}[x_i x_i^T]^{-1} \mathbb{E}[x_i y_i]$ ; the corresponding OLS estimator is  $\hat{\beta}^{\text{ols}} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$ , where  $\mathbf{X} = (x_1, \dots, x_n)^T$  is the  $n \times p$  design matrix and  $y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ .

Additionally, throughout the text we let  $[m] = \{1, 2, \dots, m\}$ , for positive integers  $m$ , and we denote the size of a set  $S$  by  $|S|$ . The  $m$ -th derivative of a function  $g : \mathbb{R} \rightarrow \mathbb{R}$  is denoted by  $g^{(m)}$ . For a vector  $u \in \mathbb{R}^p$  and a  $n \times p$  matrix  $\mathbf{U}$ , we let  $\|u\|_q$  and  $\|\mathbf{U}\|_q$  denote the  $\ell_q$ -vector and -operator norms, respectively. If  $S \subseteq [n]$ , let  $\mathbf{U}_S$  denote the  $|S| \times p$  matrix obtained from  $\mathbf{U}$  by extracting the rows that are indexed by  $S$ . For a symmetric matrix  $\mathbf{M} \in \mathbb{R}^{p \times p}$ ,  $\lambda_{\max}(\mathbf{M})$  and  $\lambda_{\min}(\mathbf{M})$  denote the maximum and minimum eigenvalues, respectively, and  $\rho_k(\mathbf{M})$  denotes the condition number of  $\mathbf{M}$  with respect to  $k$ -norm. We denote by  $\mathbf{N}_q$  the  $q$ -variate normal distribution, and

all expectations are over all randomness inside the brackets. Finally, we use  $a \lesssim b$  and  $a \leq \mathcal{O}(b)$  interchangeably, whichever is convenient (where  $\mathcal{O}(\cdot)$  refers to the big O notation).

### 3 OLS is equivalent to the true minimizer up to a scalar factor

To motivate our methodology, we assume in this section that the covariates are multivariate normal, as in [Bri82]. These distributional assumptions will be relaxed in Section 5.

**Proposition 3.1.** *Assume that the covariates are multivariate normal with mean 0 and covariance matrix  $\Sigma$ , i.e.  $x_i \sim \mathcal{N}_p(0, \Sigma)$ . Then  $\beta^{\text{pop}}$  can be written as*

$$\beta^{\text{pop}} = c_\Psi \times \beta^{\text{ols}}, \quad (3.1)$$

where  $c_\Psi \in \mathbb{R}$  is the fixed point of the mapping

$$z \rightarrow \mathbb{E} \left[ \Psi^{(2)}(\langle x_i, \beta^{\text{ols}} \rangle z) \right]^{-1}. \quad (3.2)$$

*Proof of Proposition 3.1.* The optimal point in the optimization problem (2.1), has to satisfy the following normal equations,

$$\mathbb{E}[yx_i] = \mathbb{E} \left[ x_i \Psi^{(1)}(\langle x_i, \beta \rangle) \right]. \quad (3.3)$$

Now, denote by  $\phi(x \mid \Sigma)$  the multivariate normal density with mean 0 and covariance matrix  $\Sigma$ . We recall the well-known property of Gaussian density  $d\phi(x \mid \Sigma)/dx = -\Sigma^{-1}x\phi(x \mid \Sigma)$ . Using this and integration by parts on the right hand side of the above equation, we obtain

$$\begin{aligned} \mathbb{E} \left[ x_i \Psi^{(1)}(\langle x_i, \beta \rangle) \right] &= \int x \Psi^{(1)}(\langle x, \beta \rangle) \phi(x \mid \Sigma) \, dx, \\ &= \Sigma \beta \underbrace{\mathbb{E} \left[ \Psi^{(2)}(\langle x_i, \beta \rangle) \right]}_{\in \mathbb{R}}, \end{aligned} \quad (3.4)$$

which is basically the Stein's lemma. Combining this with the normal equations (3.3) and multiplying both side with  $\Sigma^{-1}$ , we obtain the desired result.  $\square$

Proposition 3.1 and its proof provide the main intuition behind our proposed method. Observe that in our derivation, we only worked with the right hand side of the normal equations (3.3) which does not depend on the response variable  $y_i$ . Therefore, the equivalence will hold regardless of the joint distribution of  $(y_i, x_i)$ . This is the main difference from the proof of [Bri82] where  $y_i$  is assumed to follow a single index model. In Section 5, where we extend the method to non-Gaussian predictors, the identity (3.4) is generalized via the zero-bias transformations [GR97].

#### 3.1 Regularization

A version of Proposition 3.1 incorporating regularization — an important tool for datasets where  $p$  is large relative to  $n$  or the predictors are highly collinear — is also possible, as outlined briefly

---

**Algorithm 1** SLS: Scaled Least Squares Estimator

---

**Input:** Data  $(y_i, x_i)_{i=1}^n$

**Step 1. Compute the least squares estimator:**  $\hat{\beta}^{\text{ols}}$  and  $\hat{y} = \mathbf{X}\hat{\beta}^{\text{ols}}$ .

For a sub-sampling based OLS estimator, let  $S \subset [n]$  be a random subset and take  $\hat{\beta}^{\text{ols}} = \frac{|S|}{n} (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T y$ .

**Step 2. Solve the following equation for  $c \in \mathbb{R}$ :**  $1 = \frac{c}{n} \sum_{i=1}^n \Psi^{(2)}(c \hat{y}_i)$ .

Use Newton's root-finding method:

Initialize  $c$ ;

Repeat until convergence:

$$c \leftarrow c - \frac{c \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(c \hat{y}_i) - 1}{\frac{1}{n} \sum_{i=1}^n \{\Psi^{(2)}(c \hat{y}_i) + c \hat{y}_i \Psi^{(3)}(c \hat{y}_i)\}}.$$

**Output:**  $\hat{\beta}^{\text{sls}} = c \times \hat{\beta}^{\text{ols}}$ .

---

in this section. We focus on  $\ell^2$ -regularization (ridge regression) in this section; some connections with lasso ( $\ell^1$ -regularization) are discussed in Section 5 and Corollary 5.2.

For  $\lambda \geq 0$ , define the  $\ell_2$ -regularized empirical risk minimizer,

$$\beta_{\lambda}^{\text{pop}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \mathbb{E} [\Psi(\langle x_i, \beta \rangle) - y_i \langle x_i, \beta \rangle] + \frac{\lambda}{2} \|\beta\|_2^2 \quad (3.5)$$

and the corresponding  $\ell^2$ -regularized OLS coefficients  $\beta_{\lambda}^{\text{ols}} = (\mathbb{E} [x_i x_i^T] + \lambda \mathbf{I})^{-1} \mathbb{E} [x_i y_i]$  (so  $\beta^{\text{pop}} = \beta_0^{\text{pop}}$  and  $\beta^{\text{ols}} = \beta_0^{\text{ols}}$ ). The same argument as above implies that

$$\beta_{\lambda}^{\text{pop}} = c_{\Psi} \times \beta_{\gamma}^{\text{ols}}, \quad \text{where } \gamma = \lambda c_{\Psi}. \quad (3.6)$$

This suggests that the ordinary ridge regression for the linear model can be used to estimate the  $\ell^2$ -regularized empirical risk minimizer  $\beta_{\lambda}^{\text{pop}}$ . Further pursuing these ideas for problems where regularization is a critical issue may be an interesting area for future research.

## 4 SLS: Scaled Least Squares estimator

Motivated by the results in the previous section, we design a computationally efficient algorithm that approximates the stochastic optimization problem (1.1) that is as simple as solving the least squares problem; it is described in Algorithm 1. The algorithm has two basic steps. First, we estimate the OLS coefficients, and then in the second step we estimate the proportionality constant via a simple root-finding algorithm.

There are numerous fast optimization methods to solve the least squares problem, and even a superficial review of these could go beyond the page limits of this paper. We emphasize that this step (finding the OLS estimator) does not have to be iterative and it is the main computational cost of the proposed algorithm. We suggest using a sub-sampling based estimator for  $\beta^{\text{ols}}$ , where we only use a subset of the observations to estimate the covariance matrix. Let  $S \subset [n]$  be a random sub-sample and denote by  $\mathbf{X}_S$  the sub-matrix formed by the rows of  $\mathbf{X}$  in  $S$ . Then the

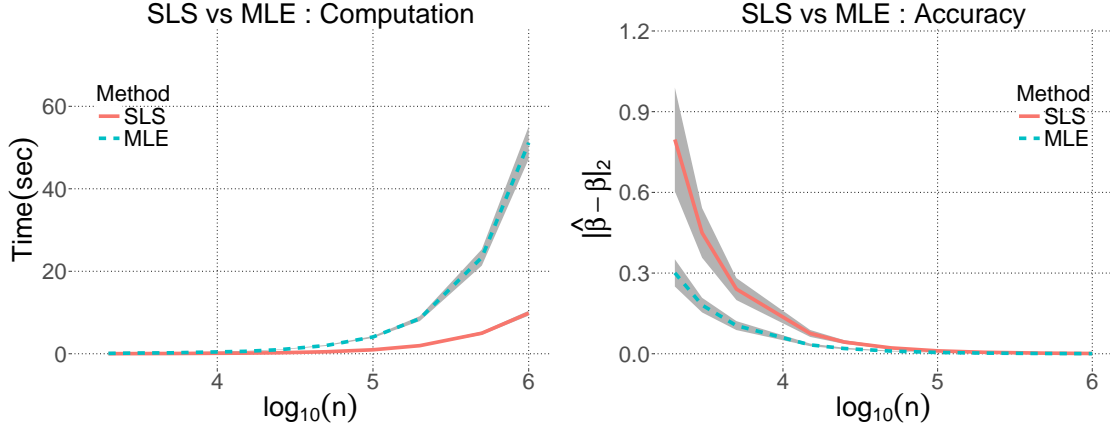


Figure 1: Logistic regression with iid standard Gaussian design. The left plot shows the computational cost (time) for finding the MLE and SLS as  $n$  grows and  $p = 200$ . The right plot depicts the accuracy of the estimators. In the regime where the MLE is expensive to compute, the SLS is found much more rapidly and has the same accuracy. R’s built-in functions are used to find the MLE.

sub-sampled OLS estimator is given as  $\hat{\beta}^{\text{ols}} = (\frac{1}{|S|} \mathbf{X}_S^T \mathbf{X}_S)^{-1} \frac{1}{n} \mathbf{X}^T y$ . Properties of sub-sampling and sketching based estimators have been well-studied [Ver10, DLFU13, EM15, PW15, RKM16]. For sub-Gaussian covariates, it suffices to use a sub-sample size of  $\mathcal{O}(p \log(p))$  [Ver10]. Hence, this step requires a single time computational cost of  $\mathcal{O}(|S|p^2 + p^3 + np) \approx \mathcal{O}(p \max\{p^2 \log(p), n\})$ . For other approaches, we refer reader to [RT08, DMMS11, DLFU13, EM15] and the references therein.

The second step of Algorithm 1 involves solving a simple root-finding problem. As with the first step of the algorithm, there are numerous methods available for completing this task. Newton’s root-finding method with quadratic convergence or Halley’s method with cubic convergence may be appropriate choices. We highlight that this step costs only  $\mathcal{O}(n)$  per-iteration and that we can attain up to a cubic rate of convergence. The resulting per-iteration cost is cheaper than other commonly used batch algorithms by at least a factor of  $\mathcal{O}(p)$  — indeed, the cost of computing the gradient is  $\mathcal{O}(np)$ . For simplicity, we use Newton’s root-finding method.

Correct initialization of the scaling constant  $c$  depends on the optimization problem. For example, in the case of GLM problems, assuming that the GLM is a good approximation to the true conditional distribution, by the law of total variance and basic properties of GLMs, we have

$$\text{Var}(y_i) = \mathbb{E}[\text{Var}(y_i | x_i)] + \text{Var}(\mathbb{E}[y_i | x_i]) \approx c_{\Psi}^{-1} + \text{Var}(\Psi^{(1)}(\langle x_i, \beta \rangle)). \quad (4.1)$$

It follows that the initialization  $c = 2/\text{Var}(y_i)$  is reasonable as long as  $c_{\Psi}^{-1} \approx \mathbb{E}[\text{Var}(y_i | x_i)]$  is not much smaller than  $\text{Var}(\Psi^{(1)}(\langle x_i, \beta \rangle))$ . Our experiments show that SLS is very robust to initialization.

In Figure 1, we compare the performance of our SLS estimator to that of the MLE in a GLM optimization problem, when both are used to analyze synthetic data generated from a logistic regression model under general Gaussian design with randomly generated covariance matrix. The left plot shows the computational cost of obtaining both estimators as  $n$  increases for fixed  $p$ . The

right plot shows the accuracy of the estimators. In the regime  $n \gg p \gg 1$  — where the MLE is hard to compute — the MLE and the SLS achieve the same accuracy, yet SLS has significantly smaller computation time. We refer the reader to Section 5 for theoretical results characterizing the finite sample behavior of the SLS.

## 5 Theoretical results

In this section, we use the zero-bias transformations [GR97] to generalize the equivalence relation given in the previous section to the settings where the covariates are non-Gaussian.

**Definition 1.** *Let  $z$  be a random variable with mean 0 and variance  $\sigma^2$ . Then, there exists a random variable  $z^*$  that satisfies  $\mathbb{E}[zf(z)] = \sigma^2 \mathbb{E}[f^{(1)}(z^*)]$ , for all differentiable functions  $f$ . The distribution of  $z^*$  is said to be the  $z$ -zero-bias distribution.*

The existence of  $z^*$  in Definition 1 is a consequence of Riesz representation theorem [GR97]. The normal distribution is the unique distribution whose zero-bias transformation is itself (i.e. the normal distribution is a fixed point of the operation mapping the distribution of  $z$  to that of  $z^*$  — which is basically Stein’s lemma).

To provide some intuition behind the usefulness of the zero-bias transformation, we refer back to the proof of Proposition 3.1. For simplicity, assume that the covariate vector  $x_i$  has iid entries with mean 0, and variance 1. Then the zero-bias transformation applied to the  $j$ -th normal equation in (3.3) yields

$$\underbrace{\mathbb{E}[y_i x_{ij}] = \mathbb{E}\left[x_{ij} \Psi^{(1)}(x_{ij} \beta_j + \sum_{k \neq j} x_{ik} \beta_k)\right]}_{j\text{-th normal equation}} = \beta_j \underbrace{\mathbb{E}\left[\Psi^{(2)}(x_{ij}^* \beta_j + \sum_{k \neq j} x_{ik} \beta_k)\right]}_{\text{Zero-bias transformation}}. \quad (5.1)$$

The distribution of  $x_{ij}^*$  is the  $x_{ij}$ -zero-bias distribution and is entirely determined by the distribution of  $x_{ij}$ ; general properties of  $x_{ij}^*$  can be found, for example, in [CGS10]. If  $\beta$  is well spread, it turns out that taken together, with  $j = 1, \dots, p$ , the far right-hand side in (5.1) behaves similar to the right side of (3.4), with  $\Sigma = \mathbf{I}$ ; that is, the behavior is similar to the Gaussian case, where the proportionality relationship given in Proposition 3.1 holds. This argument leads to an approximate proportionality relationship for problems with non-Gaussian predictors, which, when carried out rigorously, yields the following result.

**Theorem 5.1.** *Suppose that the whitened covariates  $w_i = \Sigma^{-1/2} x_i$  are independent with mean 0, covariance  $\mathbf{I}$ , and have sub-Gaussian norm bounded by  $\kappa$ . Furthermore,  $w_i$ ’s have constant first and second conditional moments, i.e.,  $\forall j \in [p]$  and  $\tilde{\beta} = \Sigma^{1/2} \beta^{\text{pop}}$ ,  $\mathbb{E}[w_{ij} | \Sigma_{k \neq j} \tilde{\beta}_k w_{ik}]$  and  $\mathbb{E}[w_{ij}^2 | \Sigma_{k \neq j} \tilde{\beta}_k w_{ik}]$  are constant. Let  $\|\beta^{\text{pop}}\|_2 = \tau$  and assume  $\beta^{\text{pop}}$  is  $r$ -well-spread in the sense that  $\tau / \|\beta^{\text{pop}}\|_\infty = r\sqrt{p}$  for some  $r \in (0, 1]$ , and the function  $\Psi^{(2)}$  is Lipschitz continuous with constant  $k$ . Then, for  $c_\Psi = 1/\mathbb{E}[\Psi^{(2)}(\langle x_i, \beta^{\text{pop}} \rangle)]$ , and  $\rho = \rho_\infty(\Sigma^{1/2})$  denoting the condition number of  $\Sigma^{1/2}$ , we have*

$$\left\| \frac{1}{c_\Psi} \times \beta^{\text{pop}} - \beta^{\text{ols}} \right\|_\infty \leq \frac{\eta}{p}, \quad \text{where } \eta = 8k\kappa^3 \rho \|\Sigma^{1/2}\|_\infty (\tau/r)^2. \quad (5.2)$$



Theorem 5.1 is proved in the Appendix. It implies that the population parameters  $\beta^{\text{ols}}$  and  $\beta^{\text{pop}}$  are approximately equivalent up to a scaling factor, with an error bound of  $\mathcal{O}(1/p)$ . The assumption that  $\beta^{\text{pop}}$  is well-spread can be relaxed with minor modifications. For example, if we have a sparse coefficient vector, where  $\text{supp}(\beta^{\text{pop}}) = \{j; \beta_j^{\text{pop}} \neq 0\}$  is the support set of  $\beta^{\text{pop}}$ , then Theorem 5.1 holds with  $p$  replaced by the size of the support set.

The assumptions on the conditional moments are the relaxed versions of assumptions that are commonly encountered in dimension reduction techniques. For example, sliced inverse regression methods assume that the first conditional moment  $\mathbb{E}[x|\langle x, \beta \rangle]$  is linear in  $x$  for all  $\beta$  [LD89, Li91], which is satisfied by elliptically distributed random vectors. An important case that is not covered by these methods is the independent coordinate case, i.e., when the whitened covariates have independent, but not necessarily identical entries. It is straightforward to observe that this case satisfies the assumptions of Theorem 5.1. We refer reader to [LD09], for a good review of dimension reduction techniques and their corresponding assumptions. We also highlight that our moment assumptions can be relaxed further, at the expense of introducing some additional complexity into the results.

An interesting consequence of Theorem 5.1 and the remarks following the theorem is that whenever an entry of  $\beta^{\text{pop}}$  is zero, the corresponding entry of  $\beta^{\text{ols}}$  has to be small, and conversely. For  $\lambda \geq 0$ , define the lasso coefficients

$$\beta_\lambda^{\text{lasso}} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \frac{1}{2} \mathbb{E}[(y_i - \langle x_i, \beta \rangle)^2] + \lambda \|\beta\|_1. \quad (5.3)$$

**Corollary 5.2.** *For any  $\lambda \geq \eta/|\text{supp}(\beta^{\text{pop}})|$ , if  $\mathbb{E}[x_i] = 0$  and  $\mathbb{E}[x_i x_i^T] = \mathbf{I}$ , we have  $\text{supp}(\beta^{\text{lasso}}) \subset \text{supp}(\beta^{\text{pop}})$ . Further, if  $\lambda$  and  $\beta^{\text{pop}}$  also satisfy that  $\forall j \in \text{supp}(\beta^{\text{pop}}), |\beta_j^{\text{pop}}| > c_\Psi (\lambda + \eta/|\text{supp}(\beta^{\text{pop}})|)$ , then we have  $\text{supp}(\beta^{\text{lasso}}) = \text{supp}(\beta^{\text{pop}})$ .*

So far in this section, we have only discussed properties of the population parameters, such as  $\beta^{\text{pop}}$  and  $\beta^{\text{ols}}$ . In the remainder of this section, we turn our attention to results for the estimators that are the main focus of this paper; these results ultimately build on our earlier results, i.e. Theorem 5.1.

In order to precisely describe the performance of  $\hat{\beta}^{\text{sls}}$ , we first need bounds on the OLS estimator. The OLS estimator has been studied extensively in the literature; however, for our purposes, we find it convenient to derive a new bound on its accuracy. While we have not seen this exact bound elsewhere, it is very similar to Theorem 5 of [DLFU13].

**Proposition 5.3.** *Assume that  $\mathbb{E}[x_i] = 0$ ,  $\mathbb{E}[x_i x_i^T] = \Sigma$ , and that  $\Sigma^{-1/2} x_i$  and  $y_i$  are sub-Gaussian with norms  $\kappa$  and  $\gamma$ , respectively. For  $\lambda_{\min}$  denoting the smallest eigenvalue of  $\Sigma$ , and  $|S| > \eta p$ ,*

$$\left\| \hat{\beta}^{\text{ols}} - \beta^{\text{ols}} \right\|_2 \leq \eta \lambda_{\min}^{-1/2} \sqrt{\frac{p}{|S|}}, \quad (5.4)$$

*with probability at least  $1 - 3e^{-p}$ , where  $\eta$  depends only on  $\gamma$  and  $\kappa$ .*

Proposition 5.3 is proved in the Supplementary Material. Our main result on the performance of  $\hat{\beta}^{\text{sls}}$  is given next.

**Theorem 5.4.** *Let the assumptions of Theorem 5.1 and Proposition 5.3 hold with  $\mathbb{E}[\|\Sigma^{-1/2}x\|_2] = \tilde{\mu}\sqrt{p}$ . Further assume that the function  $f(z) = z\mathbb{E}[\Psi^{(2)}(\langle x, \beta^{\text{ols}} \rangle z)]$  satisfies  $f(\bar{c}) > 1 + \bar{\delta}\sqrt{p}$  for some  $\bar{c}$  and  $\bar{\delta}$  such that the derivative of  $f$  in the interval  $[0, \bar{c}]$  does not change sign, i.e., its absolute value is lower bounded by  $v > 0$ . Then, for  $n$  and  $|S|$  sufficiently large, with probability at least  $1 - 5e^{-p}$ , we have*

$$\left\| \hat{\beta}^{\text{sls}} - \beta^{\text{pop}} \right\|_{\infty} \leq \eta_1 \frac{1}{p} + \eta_2 \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}}, \quad (5.5)$$

where the constants  $\eta_1$  and  $\eta_2$  are defined by

$$\eta_1 = \eta k \bar{c} \kappa^3 \rho \|\Sigma^{1/2}\|_{\infty} (\tau/r)^2 \quad (5.6)$$

$$\eta_2 = \eta \bar{c} \lambda_{\min}^{-1/2} \left( 1 + v^{-1} \lambda_{\min}^{1/2} \|\beta^{\text{ols}}\|_{\infty} \max\{(b + k/\tilde{\mu}), k\bar{c}\kappa\} \right), \quad (5.7)$$

and  $\eta > 0$  is a constant depending on  $\kappa$  and  $\gamma$ .

Note that the convergence rate of the upper bound in (5.5) depends on the sum of the two terms, both of which are functions of the data dimensions  $n$  and  $p$ . The first term on the right in (5.5) comes from Theorem 5.1, which bounds the discrepancy between  $c_{\Psi} \times \beta^{\text{ols}}$  and  $\beta^{\text{pop}}$ . This term is small when  $p$  is large, and it does not depend on the number of observations  $n$ .

The second term in the upper bound (5.5) comes from estimating  $\beta^{\text{ols}}$  and  $c_{\Psi}$ . This term is increasing in  $p$ , which reflects the fact that estimating  $\beta^{\text{pop}}$  is more challenging when  $p$  is large. As expected, this term is decreasing in  $n$  and  $|S|$ , i.e. larger sample size yields better estimates. When the full OLS solution is used ( $|S| = n$ ), the second term becomes  $\mathcal{O}(\sqrt{p \log(n)/n})$ , which suggests that  $n/\log(n)$  should be at least of order  $p$  for good performance. Also, note that there is a theoretical threshold for the sub-sampling size  $|S|$ , namely  $\mathcal{O}(n/\log(n))$ , beyond which further sub-sampling provides no improvement. This suggests that the sub-sampling size should be smaller than  $\mathcal{O}(n/\log(n))$ .

## 6 Converting One GLM to Another by Scaling

In this section, we describe an efficient algorithm to transform a generalized linear model to another. It is often the case that a practitioner would like to change the loss function (equivalently the model) he/she uses based on its performance. When the dataset is large, training a new model from the scratch is computationally inefficient and will be time consuming. In the following, we will use the proportionality relation to transition between different loss functions.

Assume that a practitioner fitted a GLM using the loss function (or cumulant generating function)  $\Psi_1$ , but he/she would like to train a new model using the loss function  $\Psi_2$ . Instead of maximizing the log-likelihood based on  $\Psi_2$ , one can exploit the proportionality relation and obtain the coefficients for the new GLM problem. Denote by  $\beta_1^{\text{pop}}$  and  $\beta_2^{\text{pop}}$  the GLM coefficients corresponding to the loss functions  $\Psi_1$  and  $\Psi_2$ , respectively. We have

$$\frac{1}{c_{\Psi_1}} \beta_1^{\text{pop}} = \frac{1}{c_{\Psi_2}} \beta_2^{\text{pop}} = \beta^{\text{ols}},$$

that is, both coefficients are proportional to the OLS coefficients which does not depend on the loss function. Therefore, these coefficients  $\beta_1^{\text{pop}}$  and  $\beta_2^{\text{pop}}$  are also proportional to each other and

---

**Algorithm 2** Conversion from one GLM to another

---

**Input:** Data  $(y_i, x_i)_{i=1}^n$ , and  $\hat{\beta}_1^{\text{glm}}$

**Step 1. Compute**  $\hat{y} = \mathbf{X}\hat{\beta}_1^{\text{glm}}$ , **and**  $\kappa = \frac{1}{n} \sum_{i=1}^n \Psi_1^{(2)}(\hat{y}_i)$ .

**Step 2. Solve the following equation for**  $\rho \in \mathbb{R}$ :  $\kappa = \frac{\rho}{n} \sum_{i=1}^n \Psi_2^{(2)}(\hat{y}_i \rho)$

Use Newton's root-finding method:

Initialize  $\rho = 1$ ;

Repeat until convergence:

$$\rho \leftarrow \rho - \frac{\rho \frac{1}{n} \sum_{i=1}^n \Psi_2^{(2)}(\rho \hat{y}_i) - \kappa}{\frac{1}{n} \sum_{i=1}^n \left\{ \Psi_2^{(2)}(\rho \hat{y}_i) + \rho \hat{y}_i \Psi_2^{(3)}(\rho \hat{y}_i) \right\}}.$$

**Output:**  $\hat{\beta}_2^{\text{glm}} = \rho \times \hat{\beta}_1^{\text{glm}}$ .

---

we can write

$$\beta_2^{\text{pop}} = \frac{c_{\Psi_2}}{c_{\Psi_1}} \beta_1^{\text{pop}} := \rho \beta_1^{\text{pop}}, \quad (6.1)$$

where the proportionality constant between two GLM types turns out to be the ratio between  $c_{\Psi_1}$  and  $c_{\Psi_2}$ , i.e.  $\rho = c_{\Psi_2}/c_{\Psi_1}$ . Using the definition of  $c_{\Psi_2}$ , we write

$$\begin{aligned} 1 &= c_{\Psi_2} \mathbb{E} \left[ \Psi_2^{(2)}(\langle x, \beta_2^{\text{pop}} \rangle) \right], \\ &= c_{\Psi_1} \rho \mathbb{E} \left[ \Psi_2^{(2)}(\langle x, \beta_1^{\text{pop}} \rangle \rho) \right]. \end{aligned}$$

Dividing the both sides by  $c_{\Psi_1}$  and using the equality  $1/c_{\Psi_1} = \mathbb{E} \left[ \Psi_1^{(2)}(\langle x, \beta_1^{\text{pop}} \rangle) \right]$ , we obtain

$$\mathbb{E} \left[ \Psi_1^{(2)}(\langle x, \beta_1^{\text{pop}} \rangle) \right] = \rho \mathbb{E} \left[ \Psi_2^{(2)}(\langle x, \beta_1^{\text{pop}} \rangle \rho) \right].$$

The above equation only involves  $\beta_1^{\text{pop}}$  as the coefficients (which is already assumed to be known or fitted by the practitioner). Therefore, if we solve it for the ratio  $\rho$ , we can estimate  $\beta_2^{\text{pop}}$  by simply using the proportionality relation given in (6.1).

The procedure described above is summarized as Algorithm 2. We emphasize that this procedure does not require the computation of the OLS estimator which was the main cost of SLS. The procedure only requires a per-iteration cost of  $\mathcal{O}(n)$ . In other words, conversion from one GLM type to another is much simpler than obtaining the GLM coefficients from the scratch.

## 7 Binary Classification with Proper Scoring Rules

In this section, we assume that for  $i \in [n]$ , the response is binary  $y_i \in \{0, 1\}$ . The binary classification problem can be described by the following minimization of an empirical risk

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i; q(\langle x_i, \beta \rangle)), \quad (7.1)$$

Table 1: Common loss functions and their canonical links

NAME	LOSS FUNCTION: $\ell(y; q)$	WEIGHT: $w(q)$	CANONICAL LINK: $q(z)$
LOG-LOSS	$-y \log(q) - (1 - y) \log(1 - q)$	$\frac{1}{q(1-q)}$	$\frac{1}{1+\exp(-z)}$
BOOSTING LOSS	$y(q^{-1} - 1)^{1/2} + (1 - y)(q^{-1} - 1)^{-1/2}$	$\frac{1}{[q(1-q)]^{3/2}}$	$\frac{1}{2} + \frac{z/2}{2(z^2/4+1)^{1/2}}$
SQUARE LOSS	$y(1 - q)^2 + (1 - y)q^2$	1	$\frac{1+z}{2}$

where  $\ell$  and  $q$  are referred to as the loss and the link functions, respectively. There are various loss functions that are used in practice. Examples include log-loss, boosting loss, square loss etc (See Table 1). As before, we constrain our analysis on the canonical links. The concept of canonical links for binary classification is introduced by [BSS05], and it is quite similar to the generalized linear problems.

For any given loss function, we define the partial losses  $\ell_k(\cdot) = \ell(y = k; \cdot)$  for  $k \in \{0, 1\}$ . Since we have a binary response variable, we can write any loss in the following format

$$\begin{aligned}\ell(y; q) &= y\ell_1(q) + (1 - y)\ell_0(q), \\ &= y(\ell_1(q) - \ell_0(q)) + \ell_0(q).\end{aligned}$$

The above formulation is of the form of a generalized linear problem. Before moving forward, we recall the concept of proper scoring in binary classification, which is sometimes referred to as Fisher consistency.

**Definition 2** (Proper scoring rules). *Assume that  $y \sim \text{Bernoulli}(\eta)$ . If the expected loss  $\mathbb{E}[\ell(y, q)]$  is minimized by  $q = \eta$  for all  $\eta \in (0, 1)$ , we call the loss function a proper scoring rule.*

The following theorem by [Sch89] provides a methodology for constructing a loss function for the proper scoring rules.

**Theorem 7.1** ([Sch89]). *Let  $w(dt)$  be a positive measure on  $(0, 1)$  that is finite on interval  $(\epsilon, 1 - \epsilon)$   $\forall \epsilon > 0$ . Then the following defines a proper scoring rule*

$$\ell_1(q) = \int_q^1 (1 - t)w(dt), \text{ and } \ell_0(q) = \int_0^q tw(dt).$$

The measure  $w(dt)$  uniquely defines the loss function (generally referred to as the weight function, since all losses can be written as weighted average of cost weighted misclassification error [BSS05, RW10]). Examples of weight functions is given in Table 1. The above theorem has many interesting interpretations; one that is most useful to us is that  $\ell_0^{(1)}(q) = qw(q)$ .

The notion of canonical links for proper scoring rules are introduced by [BSS05], which corresponds to the notion of matching loss [HKW99, RW10]. The derivation of canonical links stems from the Hessian of the above minimization, which remedies two potential problems: non-convexity and asymptotic variance inflation. It turns out that by setting  $w(q)q^{(1)}$  as constant, one

can remedy both problems [BSS05]. We will skip the derivation and, without loss of generality, assume that the canonical link-loss pair satisfies  $w(q)q^{(1)} = 1$ . Note that any loss function has a natural canonical link. The following Theorem summarizes this concept.

**Theorem 7.2** ([BSS05]). *For proper scoring rules with  $w > 0$ , there exists a canonical link function which is unique up to addition and multiplication by constants. Conversely, any link function is canonical for a unique proper scoring rule.*

The canonical link for a given loss function can be explicitly derived from the equation  $w(q)q^{(1)} = 1$ . We have provided some examples in Table 1. Using the definition of canonical link for proper scoring rules, we write the normal equations  $\frac{d}{d\beta} \mathbb{E}[\ell(y, q(\langle x, \beta \rangle))] = 0$  as

$$\begin{aligned}\mathbb{E} \left[ xq^{(1)}(\langle x, \beta \rangle) \ell_0^{(1)}(q(\langle x, \beta \rangle)) \right] &= \mathbb{E} \left[ yxq^{(1)}(\langle x, \beta \rangle) \left( \ell_0^{(1)}(q(\langle x, \beta \rangle)) - \ell_1^{(1)}(q(\langle x, \beta \rangle)) \right) \right], \\ \mathbb{E} \left[ xq^{(1)}(\langle x, \beta \rangle) q(\langle x, \beta \rangle) w(q(\langle x, \beta \rangle)) \right] &= \mathbb{E} \left[ yq^{(1)}(\langle x, \beta \rangle) w(q(\langle x, \beta \rangle)) \right], \\ \mathbb{E} [xq(\langle x, \beta \rangle)] &= \mathbb{E} [yx], \\ \Sigma \beta \mathbb{E} [q^{(1)}(\langle x, \beta \rangle)] &= \mathbb{E} [yx].\end{aligned}$$

The last equation provides us with the analog of the proportionality relation we observed in generalized linear problems. In this case, we observe that the proportionality constant becomes  $1/\mathbb{E} [q^{(1)}(\langle x, \beta \rangle)]$ . Therefore, our algorithm can be used to obtain a fast training procedure for the binary classification problems under canonical links.

## 8 Canonicalization of the Square Loss

In this section, we present a method to approximate the square loss with a canonical form. Using this canonical approximation, we can use the techniques developed in previous sections to gain computational benefits. Consider a minimization problem of the following form

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n [y_i - f(\langle x_i, \beta \rangle)]^2. \quad (8.1)$$

The above problem is commonly encountered in many machine learning tasks – specifically, in the context of neural networks, the function  $f$  is called the activation function. Here, we consider a toy example to demonstrate how our methodology can be useful in a minimization problem of the above form.

We first use Taylor series expansion around a point  $\theta$  (which should be close to  $\langle x, \beta \rangle$ ), in order to approximate the function  $f(z)$  with a linear function around  $f(\theta)$ . This way, the square loss can be approximated with a generalized linear loss. We write

$$\begin{aligned}\min_{\beta} (y - f(\langle x, \beta \rangle))^2 &= \min_{\beta} f(\langle x, \beta \rangle)^2 - 2yf(\langle x, \beta \rangle) \\ &\approx \min_{\beta} \frac{f(\langle x, \beta \rangle)^2}{2f'(\theta)} - y\langle x, \beta \rangle.\end{aligned} \quad (8.2)$$

Then, we would have

$$\Psi(z) = \frac{f(z)^2}{2f'(\theta)}, \quad (8.3)$$

and the proportionality relation given in previous sections would hold approximately. The above approximation will be accurate when the activation function is smooth around the user-specified point  $\theta$ . We suggest to use  $\theta = 0$  since when  $p$  is large and  $\beta$  is well-spread, the inner product  $\langle x, \beta \rangle$  should be close to its expectation  $\mathbb{E}[\langle x, \beta \rangle] = 0$ . This method can be used to derive proportionality relations for GLMs with non-canonical links (conditional on link being nice), and also may be of interest in non-convex optimization.

## 9 Experiments

This section contains the results of a variety of numerical studies, which show that the Scaled Least Squares estimator reaches the minimum achievable test error substantially faster than commonly used batch algorithms for finding the MLE. Both logistic and Poisson regression models (two types of GLMs) are utilized in our analyses, which are based on several synthetic and real datasets.

Below, we briefly describe the optimization algorithms for the MLE that were used in the experiments.

1. **Newton-Raphson (NR)** achieves locally quadratic convergence by scaling the gradient by the inverse of the Hessian evaluated at the current iterate. Computing the Hessian has a per-iteration cost of  $\mathcal{O}(np^2)$ , which makes it impractical for large-scale datasets.
2. **Newton-Stein (NS)** is a recently proposed second-order batch algorithm specifically designed for GLMs [Erd15a, Erd15b]. The algorithm uses Stein’s lemma and sub-sampling to efficiently estimate the Hessian with a cost of  $\mathcal{O}(np)$  per-iteration, achieving near quadratic rates.
3. **Broyden-Fletcher-Goldfarb-Shanno (BFGS)** is the most popular and stable quasi-Newton method [Nes04]. At each iteration, the gradient is scaled by a matrix that is formed by accumulating information from previous iterations and gradient computations. The convergence is locally super-linear with a per-iteration cost of  $\mathcal{O}(np)$ .
4. **Limited memory BFGS (LBFGS)** is a variant of BFGS, which uses only the recent iterates and gradients to approximate the Hessian, providing significant improvement in terms of memory usage. LBFGS has many variants; we use the formulation given in [Bis95].
5. **Gradient descent (GD)** takes a step in the opposite direction of the gradient, evaluated at the current iterate. Its performance strongly depends on the condition number of the design matrix. Under certain assumptions, the convergence is linear with  $\mathcal{O}(np)$  per-iteration cost.
6. **Accelerated gradient descent (AGD)** is a modified version of gradient descent with an additional “momentum” term [Nes83]. Its per iteration cost is  $\mathcal{O}(np)$  and its performance strongly depends on the smoothness of the objective function.

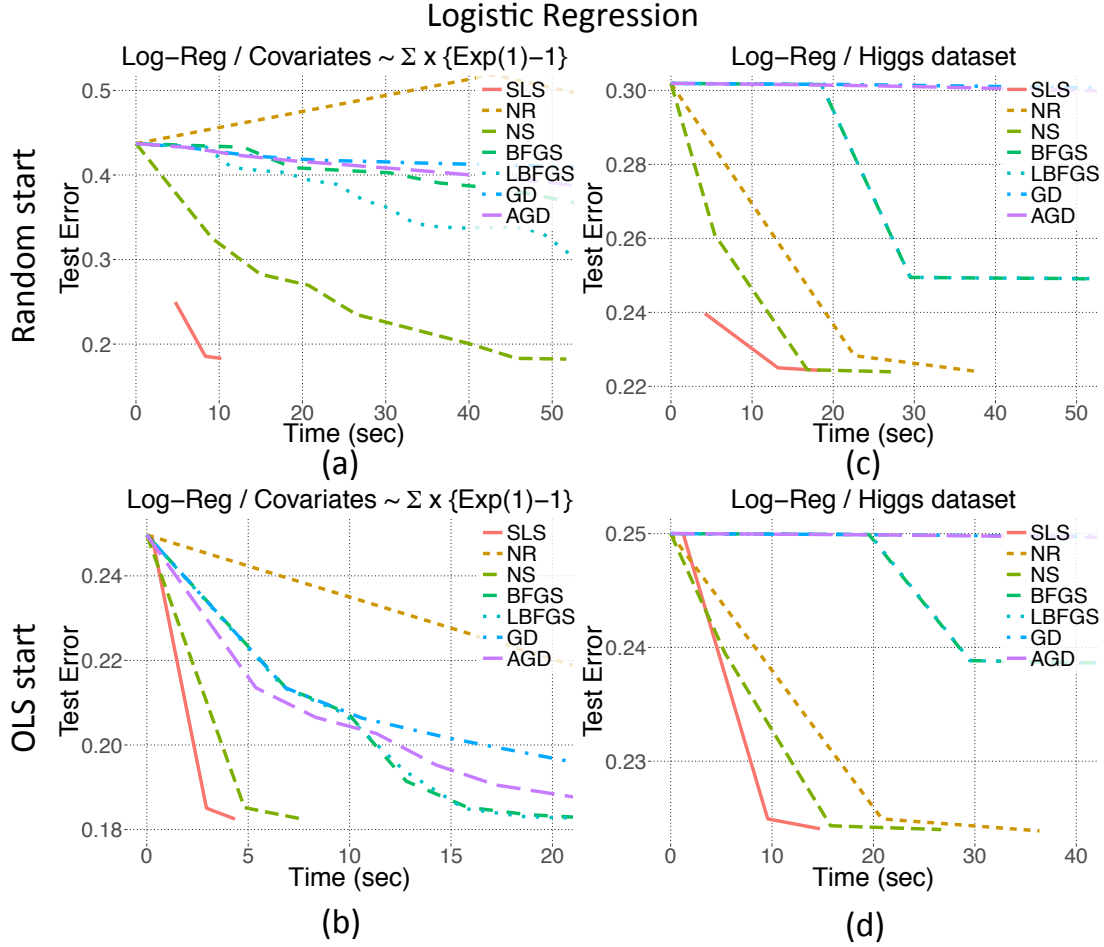


Figure 2: We compared the performance of SLS to that of MLE for the logistic regression problem on several datasets. MLE optimization is solved by various optimization algorithms. SLS is represented with red straight line. The details are provided in Table 2.

For all the algorithms for computing the MLE, the step size at each iteration is chosen via the backtracking line search [BV04].

Recall that the proposed Algorithm 1 is composed of two steps; the first finds an estimate of the OLS coefficients. This up-front computation is not needed for any of the MLE algorithms described above. On the other hand, each of the MLE algorithms requires some initial value for  $\beta$ , but no such initialization is needed to find the OLS estimator in Algorithm 1. This raises the question of how the MLE algorithms should be initialized, in order to compare them fairly with the proposed method. We consider two scenarios in our experiments: first, we use the OLS estimator computed for Algorithm 1 to initialize the MLE algorithms; second, we use a random initial value.

On each dataset, the main criterion for assessing the performance of the estimators is how rapidly the minimum test error is achieved. The test error is measured as the mean squared error

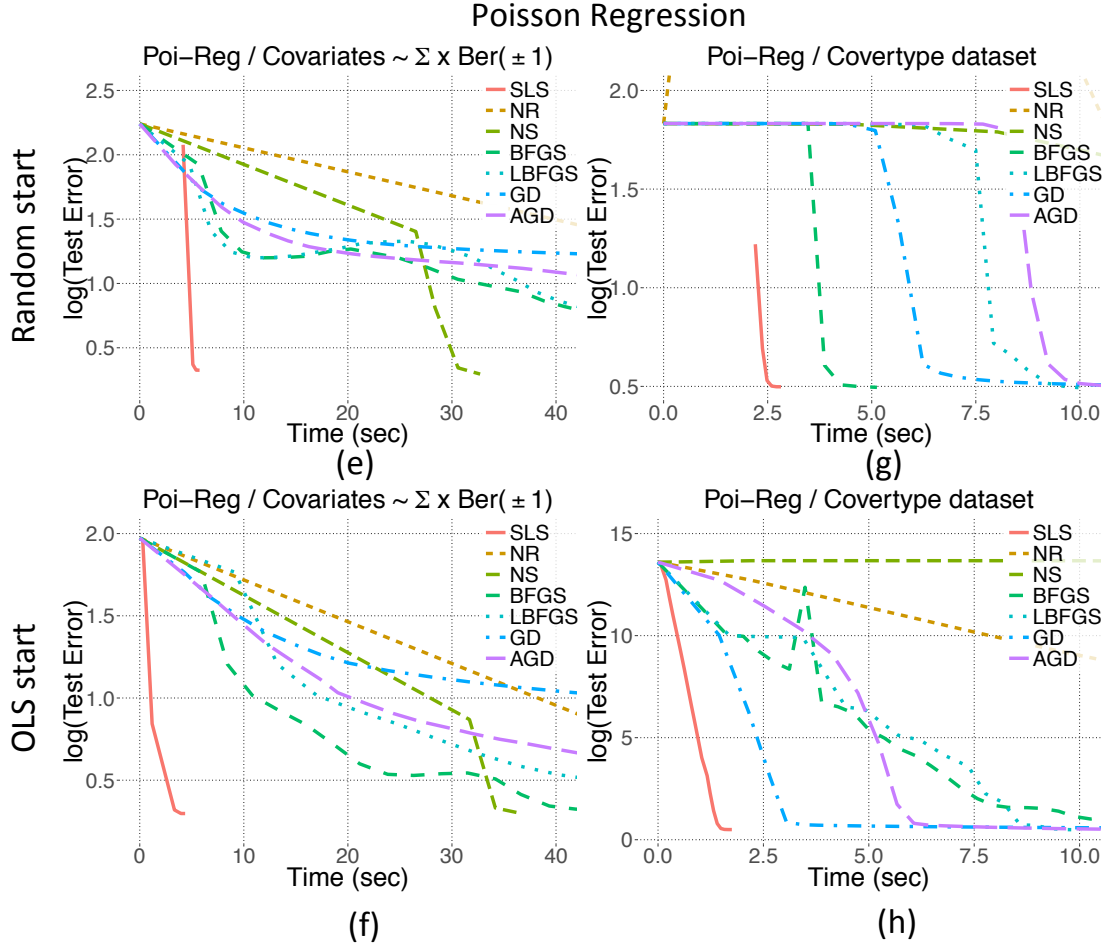


Figure 3: We compared the performance of SLS to that of MLE for the Poisson regression problem on several datasets. MLE optimization is solved by various optimization algorithms. SLS is represented with red straight line. The details are provided in Table 2.

of the estimated mean using the current parameters at each iteration on a test dataset, which is a randomly selected (and set-aside) 10% portion of the entire dataset. As noted previously, the MLE is more accurate for small  $n$  (see Figure 1). However, in the regime considered here ( $n \gg p \gg 1$ ), the MLE and the SLS perform very similarly in terms of their error rates; for instance, on the Higgs dataset, the SLS and MLE have test error rates of 22.40% and 22.38%, respectively. For each dataset, the minimum achievable test error is set to be the maximum of the final test errors, where the maximum is taken over all of the estimation methods. Let  $\Sigma^{(1)}$  and  $\Sigma^{(2)}$  be two randomly generated covariance matrices. The datasets we analyzed were: (i) a synthetic dataset generated from a logistic regression model with iid  $\{\text{exponential}(1) - 1\}$  predictors scaled by  $\Sigma^{(1)}$ ; (ii) the Higgs dataset (logistic regression) [BSW14]; (iii) a synthetic dataset generated from a Poisson regression model with iid binary( $\pm 1$ ) predictors scaled by  $\Sigma^{(2)}$ ; (iv) the Covertypes dataset (Poisson regression) [BD99].



Table 2: Details of the experiments shown in Figures 2 and 3.

MODEL	LOGISTIC REGRESSION				POISSON REGRESSION			
DATASET	$\Sigma \times \{\text{Exp}(1)-1\}$		HIGGS [BSW14]		$\Sigma \times \text{BER}(\pm 1)$		COVERTYPE [BD99]	
SIZE	$n = 6.0 \times 10^5, p = 300$		$n = 1.1 \times 10^7, p = 29$		$n = 6.0 \times 10^5, p = 300$		$n = 5.8 \times 10^5, p = 53$	
INITIALIZED	RND	OLS	RND	OLS	RND	OLS	RND	OLS
PLOT	(A)	(B)	(C)	(D)	(E)	(F)	(G)	(H)
METHOD	TIME IN SECONDS / NUMBER OF ITERATIONS (TO REACH MIN TEST ERROR)							
SLS	8.34/4	2.94/3	13.18/3	9.57/3	5.42/5	3.96/5	2.71/6	1.66/20
NR	301.06/6	82.57/3	37.77/3	36.37/3	170.28/5	130.1/4	16.7/8	32.48/18
NS	51.69/8	7.8/3	27.11/4	26.69/4	32.71/5	36.82/4	21.17/10	282.1/216
BFGS	148.43/31	24.79/8	660.92/68	701.9/68	67.24/29	72.42/26	5.12/7	22.74/59
LBFGS	125.33/39	24.61/8	6368.1/651	6946.1/670	224.6/106	357.1/88	10.01/14	10.05/17
Gd	669/138	134.91/25	100871/10101	141736/13808	1711/513	1364/374	14.35/25	33.58/87
AGD	218.1/61	35.97/12	2405.5/251	2879.69/277	103.3/51	102.74/40	11.28/15	11.95/25

In all cases, the SLS outperformed the alternative algorithms for finding the MLE by a large margin, in terms of computation. Detailed results may be found in Figures 2 and 3, and Table 2. We provide additional experiments with different datasets in the Supplementary Material.

## 10 Discussion

In this paper, we showed that the true minimizer of a generalized linear problem and the OLS estimator are approximately proportional under the general random design setting. Using this relation, we proposed a computationally efficient algorithm for large-scale problems that achieves the same accuracy as the empirical risk minimizer by first estimating the OLS coefficients and then estimating the proportionality constant through iterations that can attain quadratic or cubic convergence rate, with only  $\mathcal{O}(n)$  per-iteration cost.

We briefly mentioned that the proportionality between the coefficients holds even when there is regularization in Section 3.1. Further pursuing this idea may be interesting for large-scale problems where regularization is crucial. Another interesting line of research is to find similar proportionality relations between the parameters in other large-scale optimization problems such as support vector machines. Such relations may reduce the problem complexity significantly.

## References

- [BD99] J. A. Blackard and D. J. Dean, *Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables*, Comput. Electron. Agr. **24** (1999), 131–151.
- [Bis95] C. M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press, 1995.
- [Bri82] D. R Brillinger, *A generalized linear model with "Gaussian" regressor variables*, A Festschrift For Erich L. Lehmann, CRC Press, 1982, pp. 97–114.
- [Bro70] Charles G Broyden, *The convergence of a class of double-rank minimization algorithms 2. the new algorithm*, IMA Journal of Applied Mathematics **6** (1970), no. 3, 222–231.
- [BSS05] Andreas Buja, Werner Stuetzle, and Yi Shen, *Loss functions for binary class probability estimation and classification: Structure and applications*.
- [BSW14] Pierre Baldi, Peter Sadowski, and Daniel Whiteson, *Searching for exotic particles in high-energy physics with deep learning*, Nature communications **5** (2014).
- [BV04] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [CGS10] L. H. Y. Chen, L. Goldstein, and Q.-M. Shao, *Normal approximation by Stein's method*, Springer, 2010.
- [DLFU13] Paramveer Dhillon, Yichao Lu, Dean P Foster, and Lyle Ungar, *New subsampling algorithms for fast least squares regression*, Advances in Neural Information Processing Systems, 2013, pp. 360–368.
- [DMMS11] Petros Drineas, Michael W. Mahoney, S. Muthukrishnan, and Sarlòs, *Faster least squares approximation*, Numer. Math. **117** (2011), no. 2.
- [EBD16] Murat A Erdogdu, Mohsen Bayati, and Lee H Dicker, *Scaled least squares estimator for glms in large-scale problems*, Advances in Neural Information Processing Systems, 2016.
- [EM15] Murat A Erdogdu and Andrea Montanari, *Convergence rates of sub-sampled newton methods*, Advances in Neural Information Processing Systems, 2015, pp. 3052–3060.
- [Erd15a] Murat A Erdogdu, *Newton-stein method: A second order method for glms via stein's lemma*, Advances in Neural Information Processing Systems, 2015, pp. 1216–1224.
- [Erd15b] ———, *Newton-stein method: An optimization method for glms via stein's lemma*, arXiv preprint arXiv:1511.08895 (2015).
- [Fis36] R. A. Fisher, *The use of multiple measurements in taxonomic problems*, Annals Eugenics **7** (1936), 179–188.
- [Fle70] Roger Fletcher, *A new approach to variable metric algorithms*, The computer journal **13** (1970), no. 3, 317–322.
- [Gol70] Donald Goldfarb, *A family of variable-metric methods derived by variational means*, Mathematics of computation **24** (1970), no. 109, 23–26.
- [Gol07] L. Goldstein,  *$l^1$  bounds in normal approximation*, Annals Probability **35** (2007), 1888–1930.
- [GR97] L. Goldstein and G. Reinert, *Stein's method and the zero bias transformation with application to simple random sampling*, Annals of Applied Probability **7** (1997), 935–952.
- [HKW99] David P Helmbold, Jyrki Kivinen, and Manfred K Warmuth, *Relative loss bounds for single neurons*, IEEE Transactions on Neural Networks **10** (1999), no. 6, 1291–1304.

- [HS52] M. R. Hestenes and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, J. Res. Nat. Bur. Stand. **49** (1952), 409–436.
- [KF09] D. Koller and N. Friedman, *Probabilistic Graphical Models: Principles and Techniques*, MIT press, 2009.
- [LD89] K.-C. Li and N. Duan, *Regression analysis under link violation*, Annals of Statistics **17** (1989), 1009–1052.
- [LD09] Bing Li and Yuexiao Dong, *Dimension reduction for nonelliptically distributed predictors*, The Annals of Statistics (2009), 1272–1298.
- [Li91] Ker-Chau Li, *Sliced inverse regression for dimension reduction*, Journal of the American Statistical Association **86** (1991), no. 414, 316–327.
- [Mar10] James Martens, *Deep learning via hessian-free optimization*, Proceedings of the 27th International Conference on Machine Learning (ICML-10), 2010, pp. 735–742.
- [MN89] P. McCullagh and J. A. Nelder, *Generalized Linear Models*, 2nd ed., Chapman and Hall, 1989.
- [NB72] John A Nelder and R. Jacob Baker, *Generalized linear models*, Wiley Online Library, 1972.
- [Nes83] Y. Nesterov, *A method of solving a convex programming problem with convergence rate  $O(1/k^2)$* , Soviet Math. Dokl. **27** (1983), 372–376.
- [Nes04] ———, *Introductory Lectures on Convex Optimization: A Basic Course*, Springer, 2004.
- [PS75] C. C. Paige and M. A. Saunders, *Solution of sparse indefinite systems of linear equations*, SIAM Journal of Numerical Analysis **12** (1975), 617–629.
- [PV15] Y. Plan and R. Vershynin, *The generalized lasso with non-linear observations*, 2015, arXiv preprint arXiv:1502.04071.
- [PW15] Mert Pilanci and Martin J Wainwright, *Newton sketch: A linear-time optimization algorithm with linear-quadratic convergence*, arXiv preprint arXiv:1505.02250 (2015).
- [RKM16] Farbod Roosta-Khorasani and Michael W Mahoney, *Sub-sampled newton methods i: globally convergent algorithms*, arXiv preprint arXiv:1601.04737 (2016).
- [RT08] V. Rokhlin and M. Tygert, *A fast randomized algorithm for overdetermined linear least-squares regression*, P. Natl. Acad. Sci. **105** (2008), 13212–13217.
- [RW10] Mark D Reid and Robert C Williamson, *Composite binary losses*, Journal of Machine Learning Research **11** (2010), no. Sep, 2387–2422.
- [Sch89] Mark J Schervish, *A general method for comparing probability assessors*, The Annals of Statistics (1989), 1856–1879.
- [Sha70] David F Shanno, *Conditioning of quasi-newton methods for function minimization*, Mathematics of computation **24** (1970), no. 111, 647–656.
- [TAH15] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi, *Lasso with non-linear measurements is equivalent to one with linear measurements*, Advances in Neural Information Processing Systems, 2015, pp. 3420–3428.
- [Ver10] R. Vershynin, *Introduction to the non-asymptotic analysis of random matrices*, 2010, arXiv:1011.3027.
- [WJ08] M. J. Wainwright and M. I. Jordan, *Graphical models, exponential families, and variational inference*, Foundations and Trends in Machine Learning **1** (2008), 1–305.

## A Proof of Main Results

In this section, we provide the details and the proofs of our technical results. For convenience, we briefly state the following definitions.

**Definition 3** (Sub-Gaussian). *For a given constant  $\kappa$ , a random variable  $x \in \mathbb{R}$  is said to be sub-Gaussian if it satisfies*

$$\sup_{m \geq 1} m^{-1/2} \mathbb{E} [|x|^m]^{1/m} \leq \kappa.$$

*Smallest such  $\kappa$  is the sub-Gaussian norm of  $x$  and it is denoted by  $\|x\|_{\psi_2}$ . Similarly, a random vector  $y \in \mathbb{R}^p$  is a sub-Gaussian vector if there exists a constant  $\kappa'$  such that*

$$\sup_{v \in S^{p-1}} \|\langle y, v \rangle\|_{\psi_2} \leq \kappa'.$$

**Definition 4** (Sub-exponential). *For a given constant  $\kappa$ , a random variable  $x \in \mathbb{R}$  is called sub-exponential if it satisfies*

$$\sup_{m \geq 1} m^{-1} \mathbb{E} [|x|^m]^{1/m} \leq \kappa.$$

*Smallest such  $\kappa$  is the sub-exponential norm of  $x$  and it is denoted by  $\|x\|_{\psi_1}$ . Similarly, a random vector  $y \in \mathbb{R}^p$  is a sub-exponential vector if there exists a constant  $\kappa'$  such that*

$$\sup_{v \in S^{p-1}} \|\langle y, v \rangle\|_{\psi_1} \leq \kappa'.$$

We start with the proof of Theorem 5.1.

*Proof of Theorem 5.1.* For simplicity, we denote the whitened covariate by  $w = \Sigma^{-1/2}x$ . Since  $w$  is sub-Gaussian with norm  $\kappa$ , its  $j$ -th entry  $w_j$  has bounded third moment. That is,

$$\begin{aligned} \kappa &= \sup_{\|u\|_2=1} \|\langle u, w \rangle\|_{\psi_2}, \\ &\geq \|w_j\|_{\psi_2} = \sup_{m \geq 1} m^{-1/2} \mathbb{E} [|w_j|^m]^{1/m}, \\ &\geq \frac{1}{\sqrt{3}} \mathbb{E} [|w_j|^3]^{1/3}, \end{aligned} \tag{A.1}$$

where in the first step, we used  $u = e_j$ , the  $j$ -th standard basis vector. Hence, we obtain a bound on the third moment, i.e.,

$$\max_j \mathbb{E} [|w_j|^3] \leq 3^{3/2} \kappa^3. \tag{A.2}$$

Using the normal equations, we write

$$\begin{aligned} \mathbb{E} [yx] &= \mathbb{E} \left[ x \Psi^{(1)}(\langle x, \beta \rangle) \right] = \Sigma^{1/2} \mathbb{E} \left[ w \Psi^{(1)}(\langle w, \Sigma^{1/2} \beta \rangle) \right], \\ &= \Sigma^{1/2} \mathbb{E} \left[ w \Psi^{(1)}(\langle w, \tilde{\beta} \rangle) \right], \end{aligned} \tag{A.3}$$

where we defined  $\tilde{\beta} = \Sigma^{1/2}\beta$ . By multiplying both sides with  $\Sigma^{-1}$ , we obtain

$$\beta^{\text{ols}} = \Sigma^{-1/2} \mathbb{E} \left[ w \Psi^{(1)}(\langle w, \tilde{\beta} \rangle) \right]. \quad (\text{A.4})$$

Now we define the partial sums  $W_{-i} = \sum_{j \neq i} \tilde{\beta}_j w_j = \langle \tilde{\beta}, w \rangle - \tilde{\beta}_i w_i$ . We will focus on the  $i$ -th entry of the above expectation given in (A.4). Denoting the zero biased transformation of  $w_i$  conditioned on  $W_{-i}$  by  $w_i^*$ , we have

$$\begin{aligned} \mathbb{E} \left[ w_i \Psi^{(1)}(\langle w, \tilde{\beta} \rangle) \right] &= \mathbb{E} \left[ \mathbb{E} \left[ w_i \Psi^{(1)}(\tilde{\beta}_i w_i + W_{-i}) \mid W_{-i} \right] \right], \\ &= \tilde{\beta}_i \mathbb{E} \left[ \Psi^{(2)}(\tilde{\beta}_i w_i^* + W_{-i}) \right], \\ &= \tilde{\beta}_i \mathbb{E} \left[ \Psi^{(2)}(\tilde{\beta}_i (w_i^* - w_i) + \langle w, \tilde{\beta} \rangle) \right], \end{aligned} \quad (\text{A.5})$$

where in the second step, we used the assumption on conditional moments. Let  $\mathbf{D}$  be a diagonal matrix with diagonal entries  $\mathbf{D}_{ii} = \mathbb{E} \left[ \Psi^{(2)}(\tilde{\beta}_i (w_i^* - w_i) + \langle w, \tilde{\beta} \rangle) \right]$ . Using (A.4) together with (A.5), we obtain the equality

$$\begin{aligned} \beta^{\text{ols}} &= \Sigma^{-1/2} \mathbf{D} \tilde{\beta}, \\ &= \Sigma^{-1/2} \mathbf{D} \Sigma^{1/2} \beta. \end{aligned} \quad (\text{A.6})$$

Now, using the Lipschitz continuity assumption of the variance function, we have

$$\left| \mathbb{E} \left[ \Psi^{(2)}(\tilde{\beta}_i (w_i^* - w_i) + \langle w, \tilde{\beta} \rangle) \right] - \mathbb{E} \left[ \Psi^{(2)}(\langle w, \tilde{\beta} \rangle) \right] \right| \leq k |\tilde{\beta}_i| \mathbb{E} [|w_i^* - w_i|]. \quad (\text{A.7})$$

In the following, we will use the properties of zero-biased transformations. Consider the quantity

$$r = \sup \frac{\mathbb{E} [|w_i^* - w_i| \mid W_{-i}]}{\mathbb{E} [|w_i|^3 \mid W_{-i}]} \quad (\text{A.8})$$

where  $w_i^*$  has  $w_i$ -zero biased distribution (conditioned on  $W_{-i}$ ) and the supremum is taken with respect to all random variables with mean 0, standard deviation 1 and finite third moment, and  $w_i^*$  is achieving the minimal  $\ell_1$  coupling to  $w_i$  conditioned on  $W_{-i}$ . It is shown in [Gol07] that the above bound holds for  $r = 1.5$  for the unconditional zero-bias transformations. Here, we take a similar approach to show that the same bound holds for the conditional case as well. By using the triangle inequality, we have

$$\begin{aligned} \mathbb{E} [|w_i^* - w_i| \mid W_{-i}] &\leq \mathbb{E} [|w_i^*| \mid W_{-i}] + \mathbb{E} [|w_i| \mid W_{-i}] \\ &\leq \frac{1}{2} \mathbb{E} [|w_i|^3 \mid W_{-i}] + \mathbb{E} [|w_i|^3 \mid W_{-i}]^{1/3}. \end{aligned}$$

Since  $\mathbb{E} [|w_i|^2 \mid W_{-i}]$  is constant, it is equal to  $\mathbb{E} [|w_i|^2] = 1$ . This yields that the second term in the last line is upper bounded by  $\mathbb{E} [|w_i|^3 \mid W_{-i}]$ . Consequently, by taking expectations over both hand sides we obtain that

$$\mathbb{E} [|w_i^* - w_i|] \leq 1.5 \mathbb{E} [|w_i|^3].$$

Then the right hand side of (A.7) can be upper bounded by

$$\begin{aligned}
k|\tilde{\beta}_i|\mathbb{E}[|w_i^* - w_i|] &\leq rk \max_i \left\{ |\tilde{\beta}_i| \mathbb{E}[|w_i|^3] \right\}, \\
&\leq 1.5k \left\| \Sigma^{1/2} \beta \right\|_{\infty} 3^{3/2} \kappa^3, \\
&\leq 8k\kappa^3 \left\| \Sigma^{1/2} \beta \right\|_{\infty},
\end{aligned} \tag{A.9}$$

where in the second step we used the bound on the third moment given in (A.2). The last inequality provides us with the following result,

$$\max_i \left| \mathbf{D}_{ii} - \frac{1}{c_{\Psi}} \right| \leq 8k\kappa^3 \left\| \Sigma^{1/2} \beta \right\|_{\infty}. \tag{A.10}$$

Finally, combining this with (A.4) and (A.6), we obtain

$$\begin{aligned}
\left\| \beta^{\text{ols}} - \frac{1}{c_{\Psi}} \beta \right\|_{\infty} &= \left\| \Sigma^{-1/2} \mathbf{D} \Sigma^{1/2} \beta - \frac{1}{c_{\Psi}} \beta \right\|_{\infty}, \\
&= \left\| \Sigma^{-1/2} \left( \mathbf{D} - \frac{1}{c_{\Psi}} \mathbf{I} \right) \Sigma^{1/2} \beta \right\|_{\infty}, \\
&\leq \max_i \left| \mathbf{D}_{ii} - \frac{1}{c_{\Psi}} \right| \left\| \Sigma^{1/2} \right\|_{\infty} \left\| \Sigma^{-1/2} \right\|_{\infty} \left\| \beta \right\|_{\infty}^2, \\
&\leq 8k\kappa^3 \rho(\Sigma^{1/2}) \left\| \Sigma^{1/2} \right\|_{\infty} \frac{\tau^2}{r^2 p},
\end{aligned} \tag{A.11}$$

where in the last step, we used the assumption that  $\beta$  is  $r$ -well-spread.  $\square$

*Proof of Proposition 5.3.* For convenience, we denote the whitened covariates with  $w_i = \Sigma^{-1/2} x_i$ . We have  $\mathbb{E}[w_i] = 0$ ,  $\mathbb{E}[w_i w_i^T] = \mathbf{I}$ , and  $\|w_i\|_{\psi_2} \leq \kappa$ . Also denote the sub-sampled covariance matrix with  $\hat{\Sigma} = \frac{1}{|S|} \sum_{i \in S} x_i x_i^T$ , and its whitened version as  $\tilde{\Sigma} = \frac{1}{|S|} \sum_{i \in S} w_i w_i^T$ . Further, define  $\hat{\zeta} = \frac{1}{n} \sum_{i=1}^n w_i y_i$  and  $\zeta = \mathbb{E}[wy]$ . Then, we have

$$\hat{\beta}^{\text{ols}} = \hat{\Sigma}^{-1} \Sigma^{1/2} \hat{\zeta} \quad \text{and} \quad \beta^{\text{ols}} = \Sigma^{-1/2} \zeta.$$

For now, we work on the event that  $\hat{\Sigma}$  is invertible. We will see that this event holds with very high probability. We write

$$\begin{aligned}
\left\| \Sigma^{1/2} (\hat{\beta}^{\text{ols}} - \beta^{\text{ols}}) \right\|_2 &= \left\| \Sigma^{1/2} \hat{\Sigma}^{-1} \Sigma^{1/2} \hat{\zeta} - \Sigma^{-1/2} \zeta \right\|_2, \\
&= \left\| \tilde{\Sigma}^{-1} \left\{ \hat{\zeta} - \zeta + \left( \mathbf{I} - \Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2} \right) \zeta \right\} \right\|_2, \\
&\leq \left\| \tilde{\Sigma}^{-1} \right\|_2 \left\{ \left\| \hat{\zeta} - \zeta \right\|_2 + \left\| \mathbf{I} - \tilde{\Sigma} \right\|_2 \left\| \zeta \right\|_2 \right\},
\end{aligned} \tag{A.12}$$

where we used the triangle inequality and the properties of the operator norm.

For the first term on the right hand side of (A.12), we write

$$\begin{aligned}
\left\| \tilde{\Sigma}^{-1} \right\|_2 &= \frac{1}{\lambda_{\min}(\tilde{\Sigma})}, \\
&\leq \frac{1}{1 - \delta},
\end{aligned}$$

where we assumed that such a  $\delta > 0$  exists. In fact, when  $\delta < 0.5$ , we obtain a bound of 2 on the right hand side, which also justifies the invertibility assumption of  $\widehat{\Sigma}$ . By Lemma C.4 and the following remark, we have with probability at least  $1 - 2 \exp\{-p\}$ ,

$$\left\| \widetilde{\Sigma} - \mathbf{I} \right\|_2 \leq c \sqrt{\frac{p}{|S|}},$$

where  $c$  is a constant depending only on  $\kappa$ . When  $|S| > 4c^2p$ , we obtain

$$\left| \lambda_{\min}(\widetilde{\Sigma}) - 1 \right| \leq \left\| \widetilde{\Sigma} - \mathbf{I} \right\|_2 \leq 0.5,$$

where the first inequality follows from the Lipschitz property of the eigenvalues.

Next, we bound the difference between  $\hat{\zeta}$  and its expectation  $\zeta$ . We write the bounds on the sub-exponential norm

$$\begin{aligned} \|wy\|_{\psi_1} &= \sup_{\|v\|_2=1} \sup_{m \geq 1} m^{-1} \mathbb{E} [|\langle v, w \rangle y|^m]^{1/m}, \\ &\leq \sup_{\|v\|_2=1} \sup_{m \geq 1} m^{-1} \mathbb{E} [|\langle v, w \rangle|^{2m}]^{1/2m} \mathbb{E} [|y|^{2m}]^{1/2m}, \\ &\leq \sup_{\|v\|_2=1} \sup_{m \geq 1} m^{-1/2} \mathbb{E} [|\langle v, w \rangle|^{2m}]^{1/2m} \sup_{m \geq 1} m^{-1/2} \mathbb{E} [|y|^{2m}]^{1/2m}, \\ &\leq 2 \|w\|_{\psi_2} \|y\|_{\psi_2} = 2\gamma\kappa. \end{aligned} \tag{A.13}$$

Hence, we have  $\max_i \|w_i y_i - \mathbb{E}[w_i y_i]\|_{\psi_1} \leq 4\gamma\kappa$ . Further, let  $e_j$  denote the  $j$ -th standard basis, and notice that each entry of  $w$  is also sub-Gaussian with norm upper bounded by  $\kappa$ , i.e.,

$$\begin{aligned} \kappa &= \|w\|_{\psi_2} = \sup_{\|u\|_2=1} \|\langle u, w \rangle\|_{\psi_2}, \\ &\geq \|\langle e_j, w \rangle\|_{\psi_2} = \|w_j\|_{\psi_2}. \end{aligned} \tag{A.14}$$

Also, we can write

$$\begin{aligned} 2\gamma\kappa &\geq \|wy\|_{\psi_1} = \sup_{\|u\|_2=1} \sup_{m \geq 1} m^{-1} \mathbb{E} [|\langle u, w \rangle y|^m]^{1/m}, \\ &\geq \sup_{\|u\|_2=1} \mathbb{E} [|\langle u, w \rangle y|], \\ &\geq \sup_{\|u\|_2=1} \mathbb{E} [\langle u, w \rangle y], \\ &= \sup_{\|u\|_2=1} \langle u, \zeta \rangle = \|\zeta\|_2, \end{aligned} \tag{A.15}$$

where in the last step, we used the fact that dual norm of  $\ell_2$  norm is itself.

Next, we apply Lemma C.1 to  $\hat{\zeta} - \zeta$ , and obtain with probability at least  $1 - \exp\{-p\}$

$$\left\| \hat{\zeta} - \zeta \right\|_2 \leq c\gamma\kappa \sqrt{\frac{p}{n}},$$

whenever  $n > c^2p$  for an absolute constant  $c$ .

Combining the above results in (A.12), we obtain with probability at least  $1 - 3 \exp \{-p\}$

$$\left\| \Sigma^{1/2}(\hat{\beta}^{\text{ols}} - \beta^{\text{ols}}) \right\|_2 \leq 2 \left\{ c_1 \gamma \kappa \sqrt{\frac{p}{n}} + c_2 \gamma \kappa \sqrt{\frac{p}{|S|}} \right\} \leq \eta \sqrt{\frac{p}{|S|}} \quad (\text{A.16})$$

where  $\eta$  depends only on  $\kappa$  and  $\gamma$ , and  $|S| > \eta p$ . Finally, we write

$$\begin{aligned} \left\| \hat{\beta}^{\text{ols}} - \beta^{\text{ols}} \right\|_2 &\leq \lambda_{\min}^{-1/2} \left\| \Sigma^{1/2}(\hat{\beta}^{\text{ols}} - \beta^{\text{ols}}) \right\|_2, \\ &\leq \eta \lambda_{\min}^{-1/2} \sqrt{\frac{p}{|S|}}, \end{aligned}$$

with probability at least  $1 - 3 \exp \{-p\}$ , whenever  $|S| > \eta p$ .  $\square$

The following lemma – combined with the Proposition 5.3 – provides the necessary tools to prove Theorem 5.4.

**Lemma A.1.** *For a given function  $\Psi^{(2)}$  that is Lipschitz continuous with  $k$ , and uniformly bounded by  $b$ , we define the function  $f : \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}$  as*

$$f(c, \beta) = c \mathbb{E} \left[ \Psi^{(2)}(\langle x, \beta \rangle c) \right],$$

and its empirical counterpart as

$$\hat{f}(c, \beta) = c \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(\langle x_i, \beta \rangle c).$$

Assume that for some  $\delta, \bar{c} > 0$ , we have  $f(\bar{c}, \beta^{\text{ols}}) \geq 1 + \delta$ . Then,  $\exists c_\Psi > 0$  satisfying the equation

$$1 = f(c_\Psi, \beta^{\text{ols}}).$$

Further, assume that for some  $\tilde{\delta} > 0$ , we have  $\delta = \tilde{\delta} \sqrt{p}$ , and  $n$  and  $|S|$  sufficiently large, i.e.,

$$\min \left\{ \frac{n}{\log(n)}, |S| \right\} > K^2 / \tilde{\delta}^2$$

for  $K = \eta \bar{c} \max \{b + \kappa / \tilde{\mu}, k \bar{c} \kappa\}$ . Then, with probability  $1 - 5 \exp \{-p\}$ , there exists a constant  $\hat{c}_\Psi \in (0, \bar{c})$  satisfying the equation

$$1 = \hat{c}_\Psi \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(\langle x_i, \hat{\beta}^{\text{ols}} \rangle \hat{c}_\Psi).$$

Moreover, if the derivative of  $z \rightarrow f(z, \beta^{\text{ols}})$  is bounded below in absolute value (i.e. does not change sign) by  $v > 0$  in the interval  $z \in [0, \bar{c}]$ , then with probability  $1 - 5 \exp \{-p\}$ , we have

$$|\hat{c}_\Psi - c_\Psi| \leq C \sqrt{\frac{p}{\min \{n / \log(n), |S|\}}},$$

where  $C = K/v$ .



*Proof of Lemma A.1.* First statement is obvious. We notice that  $f(c, \beta^{\text{ols}})$  is a continuous function in its first argument with  $f(0, \beta^{\text{ols}}) = 0$  and  $f(\bar{c}, \beta^{\text{ols}}) \geq 1 + \delta$ . Hence, there exists  $c_\Psi > 0$  such that  $f(c_\Psi, \beta^{\text{ols}}) = 1$ . If there are many solutions to the above equation, we choose the one that is closest to zero. The condition on the derivative will guarantee the uniqueness of the solution.

Next, we will show the existence of  $\hat{c}_\Psi$  using a uniform concentration given by Lemma C.2. Define the ellipsoid centered around  $\beta^{\text{ols}}$  with radius  $\delta$ ,

$$\mathcal{B}_\Sigma^\delta(\beta^{\text{ols}}) = \left\{ \beta : \|\Sigma^{1/2}(\beta - \beta^{\text{ols}})\|_2 \leq \delta \right\},$$

and the event  $\mathcal{E}$  that  $\hat{\beta}^{\text{ols}}$  falls into  $\mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})$ , i.e.,

$$\mathcal{E} = \left\{ \hat{\beta}^{\text{ols}} \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}}) \right\}.$$

By Proposition 5.3 and the inequality given in (A.16), whenever  $|S| > \eta p \max\{1, \eta/\delta^2\}$ , we obtain

$$\mathbb{P}(\mathcal{E}^C) \leq 3 \exp\{-p\},$$

where  $\mathcal{E}^C$  denotes the complement of the event  $\mathcal{E}$ , and  $\eta$  is a constant depending only on  $\kappa$  and  $\gamma$ . For any  $c \in [0, \bar{c}]$ , on the event  $\mathcal{E}$ , we have

$$\left| \hat{f}(c, \hat{\beta}^{\text{ols}}) - f(c, \hat{\beta}^{\text{ols}}) \right| \leq \sup_{\beta \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})} \left| \hat{f}(c, \beta) - f(c, \beta) \right|.$$

Hence, we obtain the following inequality

$$\begin{aligned} \mathbb{P} \left( \sup_{c \in [0, \bar{c}]} \left| \hat{f}(c, \hat{\beta}^{\text{ols}}) - f(c, \hat{\beta}^{\text{ols}}) \right| > \epsilon \right) &\leq \mathbb{P} \left( \sup_{c \in [0, \bar{c}]} \left| \hat{f}(c, \hat{\beta}^{\text{ols}}) - f(c, \hat{\beta}^{\text{ols}}) \right| > \epsilon; \mathcal{E} \right) + \mathbb{P}(\mathcal{E}^C), \\ &\leq \mathbb{P} \left( \sup_{c \in [0, \bar{c}]} \sup_{\beta \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})} \left| \hat{f}(c, \beta) - f(c, \beta) \right| > \epsilon \right) + 3 \exp\{-p\}. \end{aligned}$$

In the following, we will use Lemma C.2 for the first term in the last line above. Denoting by  $w$ , the whitened covariates, we have  $\langle x, \beta \rangle = \langle w, \Sigma^{1/2} \beta \rangle$ . Therefore,

$$\begin{aligned} &\sup_{c \in [0, \bar{c}]} \sup_{\beta \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})} \left| \hat{f}(c, \beta) - f(c, \beta) \right| \\ &\leq \bar{c} \sup_{c \in [0, \bar{c}]} \sup_{\beta \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})} \left| \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(\langle w_i, \Sigma^{1/2} \beta \rangle c) - \mathbb{E} \left[ \Psi^{(2)}(\langle w, \Sigma^{1/2} \beta \rangle c) \right] \right|. \end{aligned}$$

Next, define the ball centered around  $\tilde{\beta}^{\text{ols}} = \Sigma^{1/2} \beta^{\text{ols}}$ , with radius  $\delta$  as  $\mathcal{B}_\delta(\tilde{\beta}^{\text{ols}}) = \Sigma^{1/2} \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})$ . We have  $\beta \in \mathcal{B}_\Sigma^\delta(\beta^{\text{ols}})$  if and only if  $\Sigma^{1/2} \beta \in \mathcal{B}_\delta(\tilde{\beta}^{\text{ols}})$ . Then, the right hand side of the above inequality can be written as

$$\begin{aligned} &\bar{c} \sup_{c \in [0, \bar{c}]} \sup_{\beta \in \mathcal{B}_\delta(\tilde{\beta}^{\text{ols}})} \left| \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(\langle w_i, \beta \rangle c) - \mathbb{E} \left[ \Psi^{(2)}(\langle w, \beta \rangle c) \right] \right|, \\ &= \bar{c} \sup_{\beta \in \mathcal{B}_{\bar{c}\delta}(\tilde{\beta}^{\text{ols}})} \left| \frac{1}{n} \sum_{i=1}^n \Psi^{(2)}(\langle w_i, \beta \rangle) - \mathbb{E} \left[ \Psi^{(2)}(\langle w, \beta \rangle) \right] \right|. \end{aligned}$$

Then, by Lemma C.2, we obtain

$$\mathbb{P} \left( \sup_{c \in [0, \bar{c}]} \left| \hat{f}(c, \hat{\beta}^{\text{ols}}) - f(c, \hat{\beta}^{\text{ols}}) \right| > c' \bar{c} (b + \kappa / \tilde{\mu}) \sqrt{\frac{p}{n / \log(n)}} \right) \leq 5 \exp \{-p\} \quad (\text{A.17})$$

whenever  $np > 51 \max \{\chi, \chi^{-1}\}$  where  $\chi = (b + \kappa / \tilde{\mu})^2 / (c' \delta^2 k^2 \bar{c}^2 \tilde{\mu}^2)$ .

Also, by the Lipschitz condition for  $\Psi^{(2)}$ , we have for any  $c \in [0, \bar{c}]$ , and  $\beta_1, \beta_2$ ,

$$\begin{aligned} |f(c, \beta_1) - f(c, \beta_2)| &\leq k c^2 \mathbb{E} \left[ \left| \langle w, \Sigma^{1/2}(\beta_1 - \beta_2) \rangle \right| \right] \\ &\leq k \bar{c}^2 \kappa \left\| \Sigma^{1/2}(\beta_1 - \beta_2) \right\|_2. \end{aligned}$$

Applying the above bound for  $\beta_1 = \hat{\beta}^{\text{ols}}$  and  $\beta_2 = \beta^{\text{ols}}$ , we obtain with probability  $1 - 3 \exp \{-p\}$

$$\left| f(c, \hat{\beta}^{\text{ols}}) - f(c, \beta^{\text{ols}}) \right| \leq \eta k \bar{c}^2 \kappa \sqrt{\frac{p}{|S|}}, \quad (\text{A.18})$$

where the last step follows from Proposition 5.3 and the inequality given in (A.16).

Combining this with the previous bound, and taking into account that  $\mu = \tilde{\mu} \sqrt{p}$ , for any  $c \in [0, \bar{c}]$ , with probability  $1 - 5 \exp \{-p\}$ , we obtain

$$\begin{aligned} \left| \hat{f}(c, \hat{\beta}^{\text{ols}}) - f(c, \beta^{\text{ols}}) \right| &\leq c' \bar{c} (b + \kappa / \tilde{\mu}) \sqrt{\frac{p}{n / \log(n)}} + \eta k \bar{c}^2 \kappa \sqrt{\frac{p}{|S|}} \\ &\leq K \sqrt{\frac{p}{\min \{n / \log(n), |S|\}}} \end{aligned}$$

where  $K = \eta \bar{c} \max \{b + \kappa / \tilde{\mu}, k \bar{c} \kappa\}$ . Here,  $\eta$  depends only on  $\kappa$  and  $\gamma$ .

In particular, for  $c = \bar{c}$  we observe that

$$\begin{aligned} \hat{f}(\bar{c}, \hat{\beta}^{\text{ols}}) &\geq f(\bar{c}, \beta^{\text{ols}}) - K \sqrt{\frac{p}{\min \{n / \log(n), |S|\}}} \\ &\geq 1 + \delta - K \sqrt{\frac{p}{\min \{n / \log(n), |S|\}}}. \end{aligned}$$

Therefore, for sufficiently large  $n$  and  $|S|$  satisfying

$$\min \left\{ \frac{n}{\log(n)}, |S| \right\} > K^2 / \delta^2$$

we obtain  $\hat{f}(\bar{c}, \hat{\beta}^{\text{ols}}) > 1$ . Since this function is continuous and  $\hat{f}(0, \hat{\beta}^{\text{ols}}) = 0$ , we obtain the existence of  $\hat{c}_\Psi \in [0, \bar{c}]$  with probability at least  $1 - 5 \exp \{-p\}$ .

Now, since  $\hat{c}_\Psi$  and  $c_\Psi$  satisfy the equations  $\hat{f}(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) = f(c_\Psi, \beta^{\text{ols}}) = 1$  (with high probability), by the inequality given in (A.17), with probability at least  $1 - 5 \exp \{-p\}$ , we obtain

$$\begin{aligned} \left| 1 - f(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) \right| &= \left| \hat{f}(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) - f(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) \right| \\ &\leq c' \bar{c} (b + \kappa / \tilde{\mu}) \sqrt{\frac{p}{n / \log(n)}}. \end{aligned}$$

Also, by the same argument in (A.18), and Proposition 5.3, we get

$$\begin{aligned} \left| f(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) - f(\hat{c}_\Psi, \beta^{\text{ols}}) \right| &\leq k\bar{c}^2\kappa \left\| \Sigma(\hat{\beta}^{\text{ols}} - \beta^{\text{ols}}) \right\|_2 \\ &\leq \eta k\bar{c}^2\kappa \sqrt{\frac{p}{|S|}}. \end{aligned}$$

Now, using the Taylor's series expansion of  $c \rightarrow f(c, \beta^{\text{ols}})$  around  $c_\Psi$ , and the assumption on the derivative of  $f$  with respect to its first argument, we obtain

$$\begin{aligned} v |\hat{c}_\Psi - c_\Psi| &\leq \left| f(\hat{c}_\Psi, \beta^{\text{ols}}) - f(c_\Psi, \beta^{\text{ols}}) \right| \\ &\leq \left| f(\hat{c}_\Psi, \beta^{\text{ols}}) - f(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) \right| + \left| f(\hat{c}_\Psi, \hat{\beta}^{\text{ols}}) - 1 \right| \\ &\leq \eta k\bar{c}^2\kappa \sqrt{\frac{p}{|S|}} + c'\bar{c}(b + \kappa/\tilde{\mu}) \sqrt{\frac{p}{n/\log(n)}} \\ &\leq K \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}} \end{aligned}$$

with probability at least  $1 - 5 \exp\{-p\}$ . Here, the constant  $K$  is the same as before

$$K = \eta\bar{c} \max\{b + \kappa/\tilde{\mu}, k\bar{c}\kappa\}.$$

□

*Proof of Theorem 5.4.* We have

$$\begin{aligned} \left\| \hat{\beta}^{\text{sls}} - \beta^{\text{pop}} \right\|_\infty &= \left\| \hat{c}_\Psi \hat{\beta}^{\text{ols}} - \beta^{\text{pop}} \right\|_\infty, \\ &\leq \left\| c_\Psi \beta^{\text{ols}} - \beta^{\text{pop}} \right\|_\infty + \left\| \hat{c}_\Psi \hat{\beta}^{\text{ols}} - c_\Psi \beta^{\text{ols}} \right\|_\infty, \end{aligned} \tag{A.19}$$

where we used the triangle inequality for the  $\ell_\infty$  norm. The first term on the right hand side can be bounded using Theorem 5.1. We write

$$\left\| c_\Psi \beta^{\text{ols}} - \beta^{\text{pop}} \right\|_\infty \leq \eta_1 \frac{1}{p}, \tag{A.20}$$

for  $\eta_1 = 8k\bar{c}\kappa^3\rho(\Sigma^{1/2})\|\Sigma^{1/2}\|_\infty(\tau/r)^2$ .

For the second term, we write

$$\begin{aligned} \left\| \hat{c}_\Psi \hat{\beta}^{\text{ols}} - c_\Psi \beta^{\text{ols}} \right\|_\infty &= \left\| \hat{c}_\Psi \hat{\beta}^{\text{ols}} \pm \hat{c}_\Psi \beta^{\text{ols}} - c_\Psi \beta^{\text{ols}} \right\|_\infty, \\ &\leq \left\| \hat{c}_\Psi \hat{\beta}^{\text{ols}} - \hat{c}_\Psi \beta^{\text{ols}} \right\|_\infty + \left\| \hat{c}_\Psi \beta^{\text{ols}} - c_\Psi \beta^{\text{ols}} \right\|_\infty, \\ &\leq |\hat{c}_\Psi| \left\| \hat{\beta}^{\text{ols}} - \beta^{\text{ols}} \right\|_\infty + |\hat{c}_\Psi - c_\Psi| \left\| \beta^{\text{ols}} \right\|_\infty, \end{aligned} \tag{A.21}$$

where the first step follows from triangle inequality. By Lemma A.1, for sufficiently large  $n$  and  $|S|$ , with probability  $1 - 5 \exp\{-p\}$ , the constant  $\hat{c}_\Psi$  exists and it is in the interval  $(0, \bar{c}]$ . By the same lemma, with probability  $1 - 5 \exp\{-p\}$ , we have

$$|\hat{c}_\Psi - c_\Psi| \leq \eta_4 \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}}, \tag{A.22}$$

where  $\eta_4 = \eta' v^{-1} \bar{c} \max \{b + \kappa/\tilde{\mu}, k\bar{c}\kappa\}$ , for some constant  $\eta'$  depending on the sub-Gaussian norms  $\kappa$  and  $\gamma$ .

Also, by the norm equivalence and Proposition 5.3, we have with probability  $1 - 3 \exp \{-p\}$

$$\left\| \hat{\beta}^{\text{ols}} - \beta^{\text{ols}} \right\|_{\infty} \leq \eta_3 \sqrt{\frac{p}{|S|}}, \quad (\text{A.23})$$

for  $\eta_3 = \eta'' \lambda_{\min}^{-1/2}$ , where  $\eta''$  is constant depending only on  $\gamma$  and  $\kappa$ .

Finally, combining all these inequalities with the last line of (A.19), we have with probability  $1 - 5 \exp \{-p\}$ ,

$$\begin{aligned} \left\| \hat{\beta}^{\text{sls}} - \beta^{\text{pop}} \right\|_{\infty} &\leq \eta_1 \frac{1}{p} + \eta_3 \bar{c} \sqrt{\frac{p}{|S|}} + \eta_4 \left\| \beta^{\text{ols}} \right\|_{\infty} \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}}, \\ &\leq \eta_1 \frac{1}{p} + \left( \eta_3 \bar{c} + \eta_4 \left\| \beta^{\text{ols}} \right\|_{\infty} \right) \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}}, \\ &= \eta_1 \frac{1}{p} + \eta_2 \sqrt{\frac{p}{\min\{n/\log(n), |S|\}}}, \end{aligned} \quad (\text{A.24})$$

where

$$\begin{aligned} \eta_1 &= 8k\bar{c}\kappa^3 \rho(\Sigma^{1/2}) \left\| \Sigma^{1/2} \right\|_{\infty} (\tau/r)^2 \\ \eta_2 &= \eta_3 \bar{c} + \eta_4 \left\| \beta^{\text{ols}} \right\|_{\infty}, \\ &= \eta \bar{c} \lambda_{\min}^{-1/2} \left( 1 + v^{-1} \lambda_{\min}^{1/2} \left\| \beta^{\text{ols}} \right\|_{\infty} \max \{ (b + k/\tilde{\mu}), k\bar{c}\kappa \} \right). \end{aligned} \quad (\text{A.25})$$

□

*Proof of Corollary 5.2.* The normal equations for the lasso minimization yields

$$\mathbb{E} [xx^T] \beta_{\lambda}^{\text{lasso}} - \beta^{\text{ols}} + \lambda s = 0,$$

where  $s \in \partial \left\| \beta_{\lambda}^{\text{lasso}} \right\|_1$ . It is well-known that under the orthogonal design where the covariates have i.i.d. entries, the above equation reduces to

$$\text{soft}(\beta^{\text{ols}}; \lambda) = \beta_{\lambda}^{\text{lasso}},$$

where  $\text{soft}(\cdot; \lambda)$  denotes the soft thresholding operator at level  $\lambda$ . For any  $\beta \in \mathbb{R}^p$ , let  $\text{supp}(\beta)$  denote the support of  $\beta$ , i.e., the set  $\{i \in [p] : \beta_i \neq 0\}$ . We have

$$\begin{aligned} \text{supp}(\beta_{\lambda}^{\text{lasso}}) &= \{i \in [p] : \beta_{\lambda, i}^{\text{lasso}} \neq 0\}, \\ &= \{i \in [p] : |\beta_i^{\text{ols}}| > \lambda\} \end{aligned}$$

By Theorem 5.1, we have

$$|\beta_i^{\text{ols}}| \leq \frac{1}{c_{\Psi}} |\beta_i^{\text{pop}}| + \frac{\eta}{|\text{supp}(\beta^{\text{pop}})|},$$

which implies that

$$\text{supp}(\beta_\lambda^{\text{lasso}}) \subset \left\{ i \in [p] : \frac{1}{c_\Psi} |\beta_i^{\text{pop}}| + \frac{\eta}{|\text{supp}(\beta^{\text{pop}})|} > \lambda \right\}.$$

Hence, whenever  $\lambda > \eta/|\text{supp}(\beta^{\text{pop}})|$ , we have

$$\text{supp}(\beta_\lambda^{\text{lasso}}) \subset \text{supp}(\beta^{\text{pop}}).$$

Further, we have by Theorem 5.1

$$\frac{1}{c_\Psi} |\beta_i^{\text{pop}}| \leq |\beta_i^{\text{ols}}| + \frac{\eta}{|\text{supp}(\beta^{\text{pop}})|}.$$

Hence, whenever  $|\beta_i^{\text{pop}}| > c_\Psi (\lambda + \eta/|\text{supp}(\beta^{\text{pop}})|)$ , we get  $|\beta_i^{\text{ols}}| > \lambda$ . If this condition is satisfied for any entry in the support of  $\beta^{\text{pop}}$ , the corresponding lasso coefficient will be non-zero. Therefore, we get

$$\text{supp}(\beta^{\text{pop}}) \subset \text{supp}(\beta_\lambda^{\text{lasso}})$$

under this assumption. Combining this with the previous result, we conclude the proof.  $\square$

## B Additional Experiments

In this section, we provide additional experiments. The overall setting is the same as Section 9. The only difference is that we change the sampling distribution of the datasets, which are stated in the title of each plot. As in Section 9, SLS estimator outperforms its competitors by a large margin in terms of the computation time.

The results are provided in Figures 4 and 5, and Table 3.

Table 3: Details of the experiments shown in Figures 4 and 5.

MODEL	LOGISTIC REGRESSION				POISSON REGRESSION			
DATASET	$\Sigma \times \text{BER}(\pm 1)$		$\Sigma \times \text{NORM}(0,1)$		$\Sigma \times \{\text{EXP}(1)-1\}$		$\Sigma \times \text{NORM}(0,1)$	
SIZE	$n = 6.0 \times 10^5, p = 300$		$n = 6.0 \times 10^5, p = 300$		$n = 6.0 \times 10^5, p = 300$		$n = 6.0 \times 10^5, p = 300$	
INITIALIZE	RND	OLS	RND	OLS	RND	OLS	RND	OLS
PLOT	(A)	(B)	(C)	(D)	(E)	(F)	(G)	(H)
METHOD↓	TIME IN SECONDS / NUMBER OF ITERATIONS (TO REACH MIN TEST ERROR)							
SLS	6.61/3	2.97/3	9.38/5	4.25/4	14.68/4	2.99/4	6.66/10	4.13/10
NR	222.21/6	84.08/3	186.33/6	115.76/4	218.1/6	218.9/4	364.63/9	363.4/9
NS	40.68/10	11.57/3	53.06/9	19.52/4	39.22/6	59.61/4	51.48/10	39.8/10
BFGS	125.83/33	35.41/9	155.3/48	24.78/8	46.61/20	48.71/12	92.84/36	74.22/38
LBFGS	142.09/38	44.41/12	444.62/143	21.79/7	96.53/39	50.56/12	296.4/111	228.1/117
Gd	409.9/134	79.45/22	1773.1/509	135.62/44	569.1/211	124.31/48	792.3/344	1041.1/366
AGD	177.3/159	43.76/12	359.56/95	53.73/18	157.9/57	63.16/16	74.74/32	62.21/32

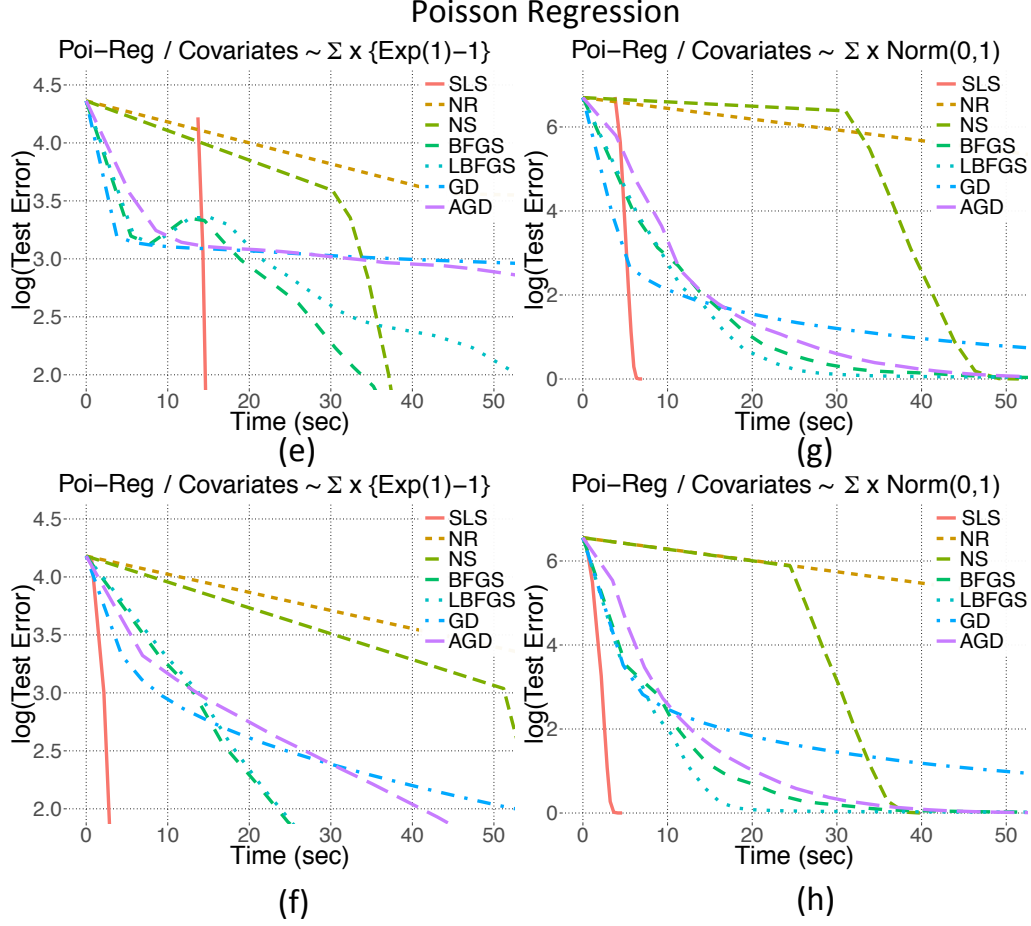


Figure 4: Additional experiments comparing the performance of SLS to that of MLE obtained with various optimization algorithms on several datasets. SLS is represented with red straight line. The details are provided in Table 3

## C Auxiliary Lemmas

**Lemma C.1** (Sub-exponential vector concentration). *Let  $x_1, x_2, \dots, x_n$  be independent centered sub-exponential random vectors with  $\max_i \|x_i\|_{\psi_1} = \kappa$ . Then we have*

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n x_i \right\|_2 > c\kappa \sqrt{\frac{p}{n}} \right) \leq \exp \{-p\}. \quad (\text{C.1})$$

whenever  $n > 4c^2p$  for an absolute constant  $c$ .

*Proof of Lemma C.1.* For a vector  $z \in \mathbb{R}^p$ , we have  $\|z\|_2 = \sup_{\|u\|_2=1} \langle u, z \rangle$  since the dual of  $\ell_2$

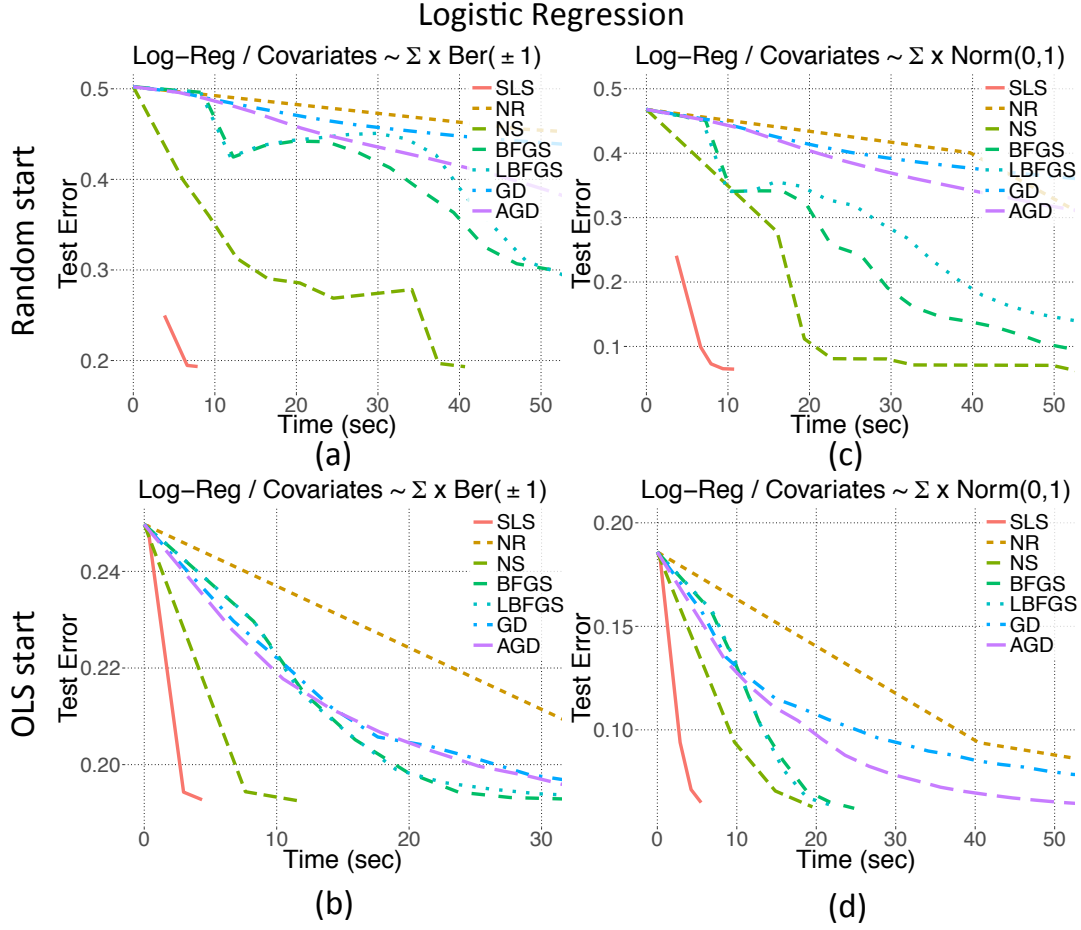


Figure 5: Additional experiments comparing the performance of SLS to that of MLE obtained with various optimization algorithms on several datasets. SLS is represented with red straight line. The details are provided in Table 3

norm is itself. Therefore, we write

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n x_i \right\|_2 > t \right) = \mathbb{P} \left( \sup_{\|u\|_2=1} \frac{1}{n} \sum_{i=1}^n \langle u, x_i \rangle > t \right).$$

Now, let  $\mathcal{N}_\epsilon$  be an  $\epsilon$ -net over  $\mathcal{S}^{p-1} = \{u \in \mathbb{R}^p : \|u\|_2 = 1\}$ , and observe that

$$\begin{aligned} \max_{u \in \mathcal{N}_\epsilon} \langle u, x \rangle &\geq (1 - \epsilon) \sup_{\|u\|_2=1} \langle u, x \rangle, \\ &= (1 - \epsilon) \|x\|_2, \end{aligned}$$

with  $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^p$ . Hence, we may write

$$\begin{aligned} \mathbb{P} \left( \sup_{\|u\|_2=1} \frac{1}{n} \sum_{i=1}^n \langle u, x_i \rangle > t \right) &\leq \mathbb{P} \left( \max_{u \in \mathcal{N}_\epsilon} \frac{1}{n} \sum_{i=1}^n \langle u, x_i \rangle > t(1 - \epsilon) \right), \\ &\leq |\mathcal{N}_\epsilon| \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n \langle u, x_i \rangle > t(1 - \epsilon) \right). \end{aligned}$$

For any  $u \in \mathcal{S}^{p-1}$ , we have  $\|\langle u, x_i \rangle\|_{\psi_1} \leq \kappa$ . Then, by the Bernstein-type inequality for sub-exponential random variables [Ver10], we have

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n \langle u, x_i \rangle > t(1 - \epsilon) \right) \leq \exp \left\{ -cn \min \left\{ \frac{t^2(1 - \epsilon)^2}{\kappa^2}, \frac{t(1 - \epsilon)}{\kappa} \right\} \right\},$$

for an absolute constant  $c$ . Therefore, the probability on the left hand side of (C.1) can be bounded by

$$\left( 1 + \frac{2}{\epsilon} \right)^p \exp \left\{ -cn \frac{t^2(1 - \epsilon)^2}{\kappa^2} \right\} = \exp \left\{ -cn \frac{t^2(1 - \epsilon)^2}{\kappa^2} + p \log \left( 1 + \frac{2}{\epsilon} \right) \right\},$$

whenever  $t < \kappa/(1 - \epsilon)$ . Choosing  $\epsilon = 0.5$  and for an absolute constant  $c' > 3.24/c$  and letting

$$t = c' \kappa \sqrt{\frac{p}{n}},$$

we conclude the proof.  $\square$

**Lemma C.2.** Let  $B(\tilde{\beta})$  denote the ball centered around  $\tilde{\beta}$  with radius  $\delta$ , i.e.,

$$B(\tilde{\beta}) = \left\{ \beta : \|\beta - \tilde{\beta}\|_2 \leq \delta \right\}.$$

For  $i = 1, \dots, n$ , let  $x_i \in \mathbb{R}^p$  be i.i.d. centered sub-Gaussian random vectors with norm bounded by  $\kappa$  and  $\mathbb{E}[\|x\|_2] = \tilde{\mu}\sqrt{p}$ . Given a function  $g : \mathbb{R} \rightarrow \mathbb{R}$  that is uniformly bounded by  $b > 0$ , and Lipschitz continuous with  $k$ ,

$$\mathbb{P} \left( \sup_{\beta \in B} \left| \frac{1}{n} \sum_{i=1}^n g(\langle x_i, \beta \rangle) - \mathbb{E}[g(\langle x, \beta \rangle)] \right| > c(b + \kappa/\tilde{\mu}) \sqrt{\frac{p}{n/\log(n)}} \right) \leq 2 \exp \{-p\},$$

whenever  $np > 51 \max\{\chi, \chi^{-1}\}$  for  $\chi = (b + \kappa/\tilde{\mu})^2/(c\delta^2 k^2 \tilde{\mu}^2)$ . Above,  $c$  is an absolute constant.

*Proof of Lemma C.2.* Let  $\mathbb{E}[\|x\|_2] = \mu = \tilde{\mu}\sqrt{p}$  and for  $\epsilon > 0$ ,  $\beta \in B(\tilde{\beta})$  and  $w \in \mathbb{R}^p$  define the bounding functions

$$\begin{aligned} l_\beta(w) &= g(\langle w, \beta \rangle) - \epsilon \|w\|_2 / 4\mu, \\ u_\beta(w) &= g(\langle w, \beta \rangle) + \epsilon \|w\|_2 / 4\mu. \end{aligned}$$

Let  $\mathcal{N}_\Delta$  be a net over  $B(\tilde{\beta})$  in the sense that for any  $\beta_1 \in B(\tilde{\beta})$ ,  $\exists \beta_2 \in \mathcal{N}_\Delta$  such that  $\|\beta_1 - \beta_2\|_2 \leq \Delta$ . We fix  $\Delta_* = \epsilon/(4k\mu)$  and write  $\forall \beta_1 \in B, \exists \beta_2 \in \mathcal{N}_{\Delta_*}$ ,



1. an upper bound of the form:

$$\begin{aligned} g(\langle w, \beta_1 \rangle) &\leq g(\langle w, \beta_2 \rangle) + k |\langle w, \beta_1 - \beta_2 \rangle|, \\ &\leq g(\langle w, \beta_2 \rangle) + k \|w\|_2 \Delta_*, \\ &= u_{\beta_2}(w), \end{aligned}$$

2. and a lower bound of the form:

$$\begin{aligned} g(\langle w, \beta_1 \rangle) &\geq g(\langle w, \beta_2 \rangle) - k |\langle w, \beta_1 - \beta_2 \rangle|, \\ &\geq g(\langle w, \beta_2 \rangle) - k \|w\|_2 \Delta_*, \\ &= l_{\beta_2}(w), \end{aligned}$$

where the second steps in the above inequalities follow from the Cauchy-Schwarz inequality. These functions are called *bracketing functions* in the context of empirical process theory.

Hence, we can write that  $\forall \beta_1 \in B(\tilde{\beta}), \exists \beta_2 \in \mathcal{N}_{\Delta_*}$  such that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n l_{\beta_2}(x_i) - \mathbb{E}[l_{\beta_2}(x)] - \epsilon/2 &\leq \frac{1}{n} \sum_{i=1}^n g(\langle x_i, \beta_1 \rangle) - \mathbb{E}[g(\langle x, \beta_1 \rangle)], \\ &\leq \frac{1}{n} \sum_{i=1}^n u_{\beta_2}(x_i) - \mathbb{E}[u_{\beta_2}(x)] + \epsilon/2. \end{aligned}$$

The above inequalities translate to the following conclusion: Whenever the following event happens,

$$\left\{ \left| \frac{1}{n} \sum_{i=1}^n g(\langle x_i, \beta_1 \rangle) - \mathbb{E}[g(\langle x, \beta_1 \rangle)] \right| > \epsilon \right\},$$

at least one of the following events happens

$$\left\{ \frac{1}{n} \sum_{i=1}^n u_{\beta_2}(x_i) - \mathbb{E}[u_{\beta_2}(x)] > \epsilon/2 \right\} \quad \text{or} \quad \left\{ \frac{1}{n} \sum_{i=1}^n l_{\beta_2}(x_i) - \mathbb{E}[l_{\beta_2}(x)] < -\epsilon/2 \right\}.$$

Therefore, using the union bound on the above events, we may obtain

$$\begin{aligned} &\mathbb{P} \left( \sup_{\beta \in B(\tilde{\beta})} \left| \frac{1}{n} \sum_{i=1}^n g(\langle x_i, \beta \rangle) - \mathbb{E}[g(\langle x, \beta \rangle)] \right| > \epsilon \right) \\ &\leq \mathbb{P} \left( \max_{\beta \in \mathcal{N}_{\Delta_*}} \frac{1}{n} \sum_{i=1}^n u_{\beta}(x_i) - \mathbb{E}[u_{\beta}(x)] > \epsilon/2 \right) \\ &\quad + \mathbb{P} \left( \max_{\beta \in \mathcal{N}_{\Delta_*}} \frac{1}{n} \sum_{i=1}^n l_{\beta}(x_i) - \mathbb{E}[l_{\beta}(x)] < -\epsilon/2 \right). \end{aligned} \tag{C.2}$$

Note that the right hand side of the above inequality has two terms both of which are of the same form. For simplicity, we bound only the first one. The bound for the second one follows from the exact same steps.

The relation between sub-Gaussian and sub-exponential norms [Ver10] allows us to write

$$\begin{aligned} \|\|x\|_2\|_{\psi_2}^2 &\leq \|\|x\|_2^2\|_{\psi_1} \leq \sum_{i=1}^p \|x_i^2\|_{\psi_1}, \\ &\leq 2 \sum_{i=1}^p \|x_i\|_{\psi_2}^2 \leq 2\kappa^2 p, \end{aligned} \tag{C.3}$$

where the second step follows from the triangle inequality. Hence, we conclude that  $\|x\|_2 - \mathbb{E}[\|x\|_2]$  is a centered sub-Gaussian random variable with norm upper bounded by  $3\kappa\sqrt{p}$ .

For  $\epsilon < 4/3$ , we notice that the random variable  $u_\beta(x) = g(\langle x, \beta \rangle) + \epsilon\|x\|_2/4\mu$  is also sub-Gaussian with norm

$$\begin{aligned} \|u_\beta(x)\|_{\psi_2} &\leq b + \frac{\epsilon}{4\tilde{\mu}} 3\kappa \\ &\leq b + \kappa/\tilde{\mu}, \end{aligned}$$

and consequently, the centered random variable  $u_\beta(x) - \mathbb{E}[u_\beta(x)]$  has the sub-Gaussian norm upper bounded by  $2b + 2\kappa/\tilde{\mu}$ .

Then, by the Hoeffding-type inequality for the sub-Gaussian random variables, we obtain

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n u_\beta(x_i) - \mathbb{E}[u_\beta(x)] > \epsilon/2\right) \leq \exp\left\{-cn \frac{\epsilon^2}{(b + \kappa/\tilde{\mu})^2}\right\}$$

for an absolute constant  $c > 0$ .

By the same argument above, one can obtain the same result for the function  $l_\beta(x)$ . Using Hoeffding bounds in (C.2) along with the union bound over the net, we immediately obtain

$$\mathbb{P}\left(\sup_{\beta \in B(\tilde{\beta})} \left|\frac{1}{n} \sum_{i=1}^n g(\langle x_i, \beta \rangle) - \mathbb{E}[g(\langle x, \beta \rangle)]\right| > \epsilon\right) \leq 2|\mathcal{N}_{\Delta_*}| \exp\left\{-cn \frac{\epsilon^2}{(b + \kappa/\tilde{\mu})^2}\right\}$$

for some absolute constant  $c$ .

Using a standard covering argument over the net  $\mathcal{N}_{\Delta_*}$  as given in Lemma C.3, we have

$$|\mathcal{N}_{\Delta_*}| \leq \left(\frac{\delta\sqrt{p}}{\Delta_*}\right)^p = \left(\frac{4\delta k\tilde{\mu}p}{\epsilon}\right)^p.$$

Combining this with the previous bound, and choosing

$$\epsilon^2 = \frac{p}{n} \frac{(b + \kappa/\tilde{\mu})^2}{2c} \log\left(\frac{32c\delta^2 k^2 \tilde{\mu}^2 p n}{(b + \kappa/\tilde{\mu})^2}\right)$$

we get

$$\begin{aligned} &2 \left(\frac{4\delta k\tilde{\mu}p}{\epsilon}\right)^p \exp\left\{-cn \frac{\epsilon^2}{(b + \kappa/\tilde{\mu})^2}\right\} \\ &= 2 \exp\left\{-\frac{p}{2} \log \log\left(\frac{32c\delta^2 k^2 \tilde{\mu}^2 p n}{(b + \kappa/\tilde{\mu})^2}\right)\right\} \\ &\leq 2 \exp\{-p\}, \end{aligned}$$

whenever  $np > 51 \max\{\chi, \chi^{-1}\}$  for  $\chi = (b + \kappa/\tilde{\mu})^2/(c\delta^2 k^2 \tilde{\mu}^2)$ .

□

**Lemma C.3** ([EM15]). *Let  $B \subset \mathbb{R}^p$  be the ball of radius  $\delta$  centered around some  $\beta \in \mathbb{R}^p$  and  $\mathcal{N}_\epsilon$  be an  $\epsilon$ -net over  $B$ . Then,*

$$|\mathcal{N}_\epsilon| \leq \left( \frac{\delta\sqrt{p}}{\epsilon} \right)^p.$$

*Proof of Lemma C.3.* The set  $B$  can be contained in a  $p$ -dimensional cube of size  $2\delta$ . Consider a grid over this cube with mesh width  $2\epsilon/\sqrt{p}$ . Then  $B$  can be covered with at most  $(2\delta/(2\epsilon/\sqrt{p}))^p$  many cubes of edge length  $2\epsilon/\sqrt{p}$ . If one takes the projection of the centers of such cubes onto  $B$  and considers the circumscribed balls of radius  $\epsilon$ , we may conclude that  $B$  can be covered with at most

$$\left( \frac{2\delta}{2\epsilon/\sqrt{p}} \right)^p$$

many balls of radius  $\epsilon$ . □

**Lemma C.4** (Corollary 5.50 of [Ver10]). *Let  $w_1, w_2, \dots, w_n$  be isotropic random vectors with sub-Gaussian norm upper bounded by  $\kappa$ . Then for every  $t > 0$ , with probability at least  $1 - 2\exp\{-c_1 t^2\}$ , the empirical covariance  $\tilde{\Sigma}$  satisfies,*

$$\left\| \tilde{\Sigma} - \mathbf{I} \right\|_2 \leq \max\{\delta, \delta^2\} \quad \text{where} \quad \delta = c_2 \sqrt{\frac{p}{n}} + \frac{t}{\sqrt{n}}$$

where  $c_1, c_2$  are constants depending only on  $\kappa$ .

**Remark 1.** For  $t = \sqrt{p/c_1}$ , we get with probability at least  $1 - 2\exp\{-p\}$ ,

$$\left\| \tilde{\Sigma} - \mathbf{I} \right\|_2 \leq C \sqrt{\frac{p}{n}}$$

where

$$C = \left\{ c_2 + \frac{1}{\sqrt{c_1}} \right\},$$

and  $n > C^2 p$ . Here,  $C$  only depends on  $\kappa$ .

**Lemma C.5** (Corollary 5.52 of [Ver10]). *Let  $x_1, x_2, \dots, x_n$  be random vectors with mean 0 and covariance  $\Sigma$  supported on a centered Euclidean ball of radius  $\sqrt{R}$ , i.e.,  $\|x_i\|_2 \leq \sqrt{R}$ . For  $\epsilon \in (0, 1)$  and  $c > 0$  an absolute constant, with probability at least  $1 - 1/p^2$ , the empirical covariance matrix satisfies*

$$\left\| \hat{\Sigma} - \Sigma \right\|_2 \leq \epsilon \|\Sigma\|_2,$$

for  $n > cR \log(p)/(\epsilon^2 \|\Sigma\|_2)$ .